



# KLEENE MEETS CHURCH: REGULAR EXPRESSIONS AS TYPES

BJØRN BUGGE GRATHWOHL, ULRIK TERP RASMUSSEN, FRITZ HENGLIN

{bugge,dolle,henglein}@diku.dk <http://www.diku.dk/kmc>

DEPARTMENT OF  
COMPUTER SCIENCE  
UNIVERSITY OF COPENHAGEN

## REGULAR EXPRESSION USAGE

Regular expressions are amongst the most widely used DSLs along with Excel and SQL. Despite this, they are ill-behaved from a programmer's perspective. The language interpretation of regular expressions as regular languages makes it difficult to use REs as a *data extraction* tool, even though this is what they are often used for.

For example, the expression  $a^*b + (a + b)^*$  is not the same as  $(a + b)^*$ , but they denote the same language

$$\mathcal{L}[a^*b + (a + b)^*] = \mathcal{L}[(a + b)^*].$$

In Perl-style implementations, this problem is addressed with *subgroup matching*:

$$\underbrace{(a^*b)}_1 \mid \underbrace{(a \mid b)^*}_2$$

but this behaves poorly in conjunction with Kleene stars:

$$\underbrace{abababab}_1 \underbrace{b}_2 \underbrace{?}_1$$

lost!

## TYPE INTERPRETATION

An RE  $E$  can be interpreted as a *type*, with  $\mathcal{T}[E]$  denoting a set of *parse trees*:

$$\begin{aligned} \mathcal{T}[0] &= \emptyset \\ \mathcal{T}[1] &= \{()\} \\ \mathcal{T}[a] &= \{a\} \\ \mathcal{T}[E_1 E_2] &= \mathcal{T}[E_1] \times \mathcal{T}[E_2] \\ \mathcal{T}[E_1 + E_2] &= \mathcal{T}[E_1] + \mathcal{T}[E_2] \\ \mathcal{T}[E_1^*] &= \{[v_1, \dots, v_n] \mid v_i \in \mathcal{T}[E_1]\} \end{aligned}$$

A parser tells the programmer *how* the input string matches  $E$ ; with  $a^*b + (a + b)^*$  and "aaaaab" it is one of:

$$\text{inl } \langle [a, a, a, a, a], b \rangle$$

$$\text{inr } [\text{inl } a, \text{inl } a, \text{inl } a, \text{inl } a, \text{inl } a, \text{inr } b].$$

Define bit-codes  $\mathcal{B}[\cdot]$  for typed trees such that:

$$\begin{aligned} \forall E. \mathcal{T}[E] &\cong \mathcal{B}[E] \\ &\subseteq \{0, 1\}^* \end{aligned}$$

## AUTOMATA REPRESENTATION

**Theorem 1** ([2, 3]). *Paths in a Thompson NFA correspond 1-to-1 to parse trees.*

**Theorem 2** ([3]). *The greedy parse tree corresponds to the "leftmost" path.*

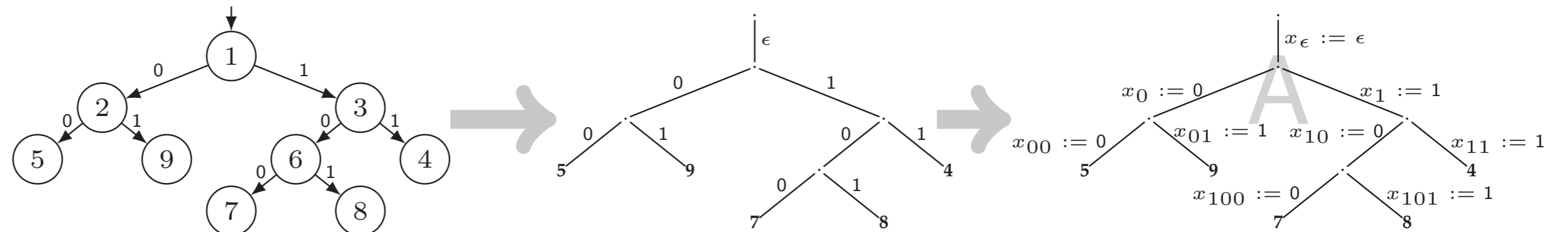
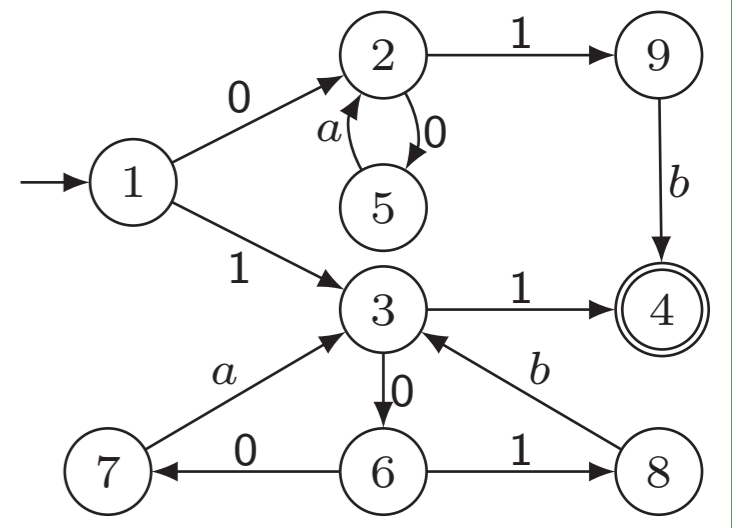
So: *bit-code of greedy parse*  $\longleftrightarrow$  *leftmost path*. We can annotate a Thompson NFA to output bit-codes, producing an *aNFA*.

## REFERENCES

- [1] Alur, R., Černý, P.: Expressiveness of Streaming String Transducers. *LIPICs*, pp. 1–12. 2010.
- [2] Brüggemann-Klein, A.: Regular Expressions into Finite Automata. *Theor. Comput. Sci.* 120(2), pp. 197–213. 1993.
- [3] Grathwohl, N.B.B., Henglein, F., Nielsen, L., Rasmussen, U.T.: Two-Pass Greedy Regular Expression Parsing. *Proc. CIAA 2013, LNCS*, vol. 7982, pp. 60–71. July 2013.
- [4] Grathwohl, N.B.B., Henglein, F., Rasmussen, U.T.: Optimally Streaming Greedy Regular Expression Parsing. *Proc. ICTAC 2014, LNCS*, vol. 8687, pp. 224–240. September 2014.
- [5] Henglein, F., Nielsen, L.: Regular Expression Containment: Coinductive Axiomatization and Computational Interpretation. *Proc. POPL 2011*, pp. 385–398. January 2011.

## PATH TREES

We maintain the set of paths through an aNFA in a *path tree*. Internal nodes represent states wherefrom paths differ. If we associate a buffer with each edge, it suffices to store the structure of the tree and its leaf nodes, not the internal nodes. As the internal nodes represent the points where paths differ, the contents of the root buffer prefixes all parse trees that can be produced by reading the remaining input. Hence, it can be output immediately.

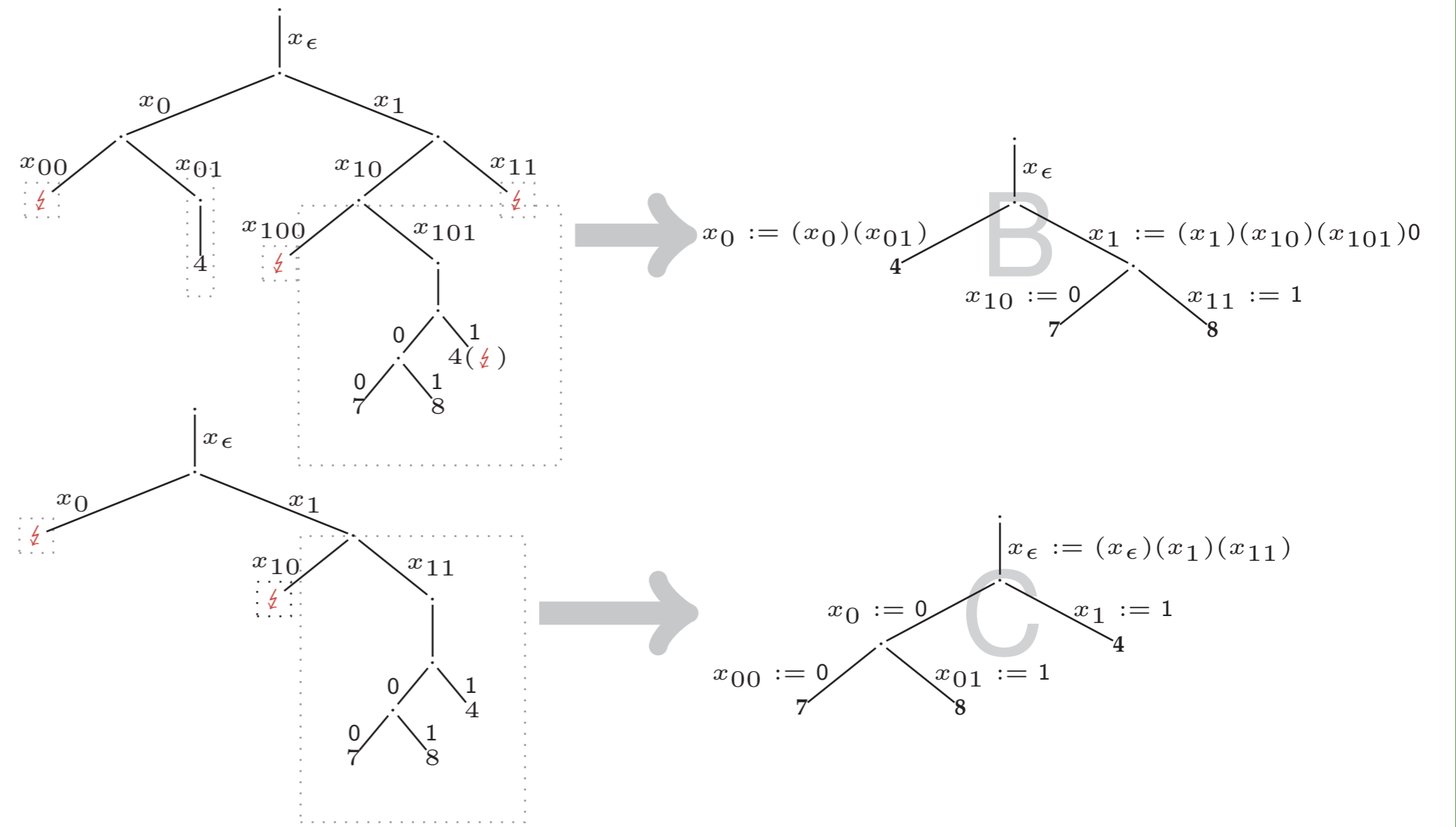


Abstract an  $\epsilon$ -closure and leave only the labelled tree structure and the leaves.

Associate a buffer to each internal edge of the path tree.

## SIMULATION

For each input symbol, the path tree is expanded at the leaves by following the aNFA transitions. We maintain uniqueness of the leaves by marking repetitions as dead; this corresponds to the greedy-leftmost disambiguation strategy. After a tree has been expanded it is contracted by combining buffers on the paths from the root to the leaves.



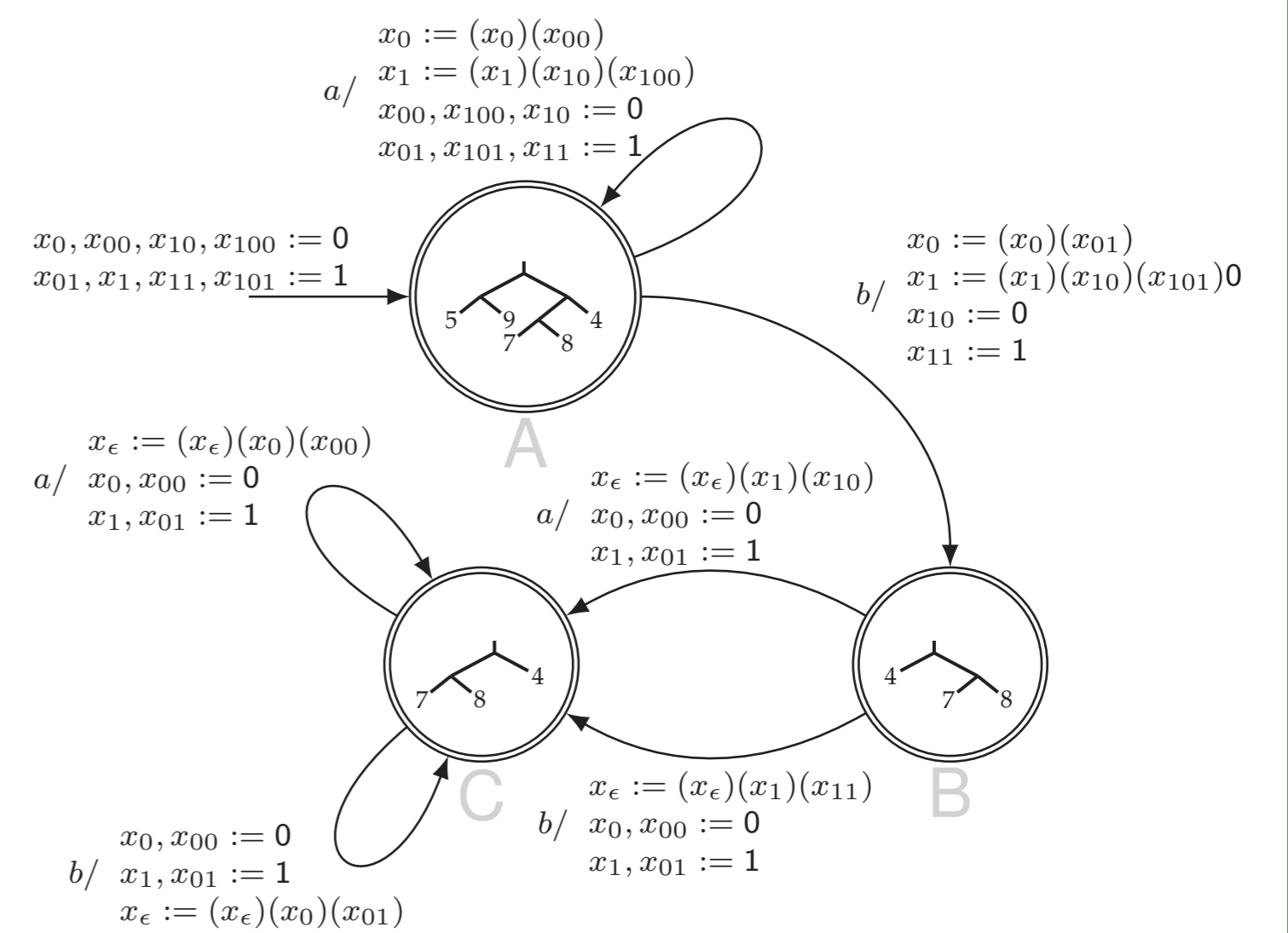
After consuming one "b" some leaves are dead, marked  $\downarrow$ . The tree is then contracted.

After consuming two "b"s and contraction.

## DETERMINIZATION

We can use the path tree simulation to determinize and disambiguate an aNFA, using the path trees as states and associating buffer updates with each transition. This corresponds to the streaming string transducers of Alur et al. [1]. The root buffer  $x_\epsilon$  represents *output*. Such a transducer thereby does *streaming parsing* of regular expressions.

The pictured transducer outputs something immediately after the first "b" has been read and the next input symbol is known. Only then can we know which of  $a^*b$  and  $(a + b)^*$  parses the input.



## OPTIMAL STREAMING

Let  $C_E(w)$  be the set of strings prefixed by  $w$  that is in the language of  $E$ . The *optimally* streaming function can be formulated as follows:

**Definition 1** (Optimally streaming function [4]). *The optimally streaming function corresponding to  $P_E(\cdot)$  is*

$$O_E(w) = \bigcap \{P_E(w') \mid w' \in C_E(w)\}$$

Intuitively, at reading input  $w$ , an optimally streaming parsing function must output the *longest* prefix of the set of paths reaching an accepting state. It requires a PSPACE-complete analysis to get optimal streaming, but for non-pathological expressions our approach suffices.