

The Average Transition Complexity of Glushkov and Partial Derivative Automata ^{*}

Sabine Broda, António Machiavelo, Nelma Moreira, Rogério Reis
sbb@dcc.fc.up.pt, ajmachia@fc.up.pt, {nam,rvr}@dcc.fc.up.pt

CMUP, Faculdade de Ciências da Universidade do Porto
Rua do Campo Alegre, 4169-007 Porto, Portugal

Abstract. In this paper, the relation between the Glushkov automaton (\mathcal{A}_{pos}) and the partial derivative automaton (\mathcal{A}_{pd}) of a given regular expression, in terms of transition complexity, is studied. The average transition complexity of \mathcal{A}_{pos} was proved by Nicaud to be linear in the size of the corresponding expression. This result was obtained using an upper bound of the number of transitions of \mathcal{A}_{pos} . Here we present a new quadratic construction of \mathcal{A}_{pos} that leads to a more elegant and straightforward implementation, and that allows the exact counting of the number of transitions. Based on that, a better estimation of the average size is presented. Asymptotically, and as the alphabet size grows, the number of transitions per state is on average 2.

Broda *et al.* computed an upper bound for the ratio of the number of states of \mathcal{A}_{pd} to the number of states of \mathcal{A}_{pos} , which is about $\frac{1}{2}$ for large alphabet sizes. Here we show how to obtain an upper bound for the number of transitions in \mathcal{A}_{pd} , which we then use to get an average case approximation. Some experimental results are presented that illustrate the quality of our estimate.

1 Introduction

The conversion methods of regular expressions into equivalent nondeterministic finite automata (NFA) are normally divided in two classes depending on whether ε -transitions are allowed or not in the resulting NFA. Paradigmatic methods of each class are the Thompson's and Glushkov's constructions, respectively. Several optimizations and worst-case descriptive and computational complexity results were obtained for both methods (see Holzer and Kutrib [HK10], and the works cited therein). Given a regular expression with n letters the size of a ε -NFA can be, in the worst-case, $\Theta(n)$. While the size of a Glushkov automaton can be $\Theta(n^2)$, $\Omega(n \log n^2)$ was proved to be a lower bound for the size of a ε -free NFA. In this context, and for practical purposes, it is useful to carry out average-case analysis, both for descriptive and computational complexities, of these methods.

^{*} This work was partially funded by Fundação para a Ciência e Tecnologia (FCT) and Program POSI, and project CANTE (PTDC/EIA-CCO/101904/2008).

The framework of analytic combinatorics, by relating the enumeration of combinatorial objects to the algebraic and complex analytic properties of generating functions, provides a powerful tool for asymptotic average-case analysis. Using this framework, Nicaud [Nic09] proved that the average transition complexity of the Glushkov automaton (\mathcal{A}_{pos}) of a regular expression α of size n is $\Theta(n)$. This result was obtained using an upper bound of the number of transitions of \mathcal{A}_{pos} . Here we present a new quadratic construction of the \mathcal{A}_{pos} that leads to a more elegant and straightforward implementation, and that allows the exact counting of the number of transitions. Based on that, a better estimation of the average size is presented. Asymptotically, and as the alphabet size grows, the number of transitions per state is on average 2.

The partial derivative automaton (\mathcal{A}_{pd}) is a quotient of the \mathcal{A}_{pos} , and thus the states of the former can be seen as mergings of states of the latter. In a previous paper [BMMR11b], we presented a technique for estimating some of those state mergings. This enabled us, in the framework of analytic combinatorics, to compute an upper bound for the ratio of the number of states of \mathcal{A}_{pd} to the number of states of \mathcal{A}_{pos} , which is about $\frac{1}{2}$ for large alphabet sizes. This upper bound was obtained by estimating the number of regular expressions that have ε as a partial derivative. In this paper, we use an analogous approach to compute an upper bound for the number of transitions in \mathcal{A}_{pd} , and study its asymptotic behaviour. As the alphabet size grows, this upper bound tends to the number of letters of the regular expression, thus it is half the number of transitions in \mathcal{A}_{pos} . The comparison with some experimental results suggests that this upper bound has an error of less than 15%.

2 Regular Expressions and Automata

In this section we briefly review some basic definitions about regular expressions and finite automata. For more details, we refer the reader to Kozen [Koz97] or Sakarovitch [Sak09].

Given an alphabet (set of *letters*) $\Sigma = \{\sigma_1, \dots, \sigma_k\}$ of size k , the set \mathcal{R} of *regular expressions*, α , over Σ is defined by the following grammar:

$$\alpha := \emptyset \mid \varepsilon \mid \sigma_1 \mid \dots \mid \sigma_k \mid (\alpha + \alpha) \mid (\alpha \cdot \alpha) \mid \alpha^* \quad (1)$$

where the operator \cdot (concatenation) is often omitted. The language associated to α is denoted by $\mathcal{L}(\alpha)$ and defined as usual. The *size* $|\alpha|$ of $\alpha \in \mathcal{R}$ is the number of symbols in α (parentheses not counted); the *alphabetic size* $|\alpha|_{\Sigma}$ is its number of letters. For example, for $\tau = (a + b)(a^* + ba^* + b^*)^*$ one has $|\tau| = 15$ and $|\tau|_{\Sigma} = 6$. We define $\varepsilon(\alpha)$ by $\varepsilon(\alpha) = \varepsilon$ if $\varepsilon \in \mathcal{L}(\alpha)$, and $\varepsilon(\alpha) = \emptyset$ otherwise. Also, we denote by α_{ε} and $\alpha_{\bar{\varepsilon}}$, respectively, the regular expressions such that $\varepsilon(\alpha_{\varepsilon}) = \varepsilon$ and $\varepsilon(\alpha_{\bar{\varepsilon}}) = \emptyset$.

A *non-deterministic automaton* (NFA) \mathcal{A} is a quintuple $(Q, \Sigma, \delta, q_0, F)$, where Q is a finite set of states, Σ is the alphabet, $\delta \subseteq Q \times \Sigma \times Q$ the transition relation, q_0 the initial state, and $F \subseteq Q$ the set of final states. The *size* of a NFA is $|Q| + |\delta|$. For $q \in Q$ and $\sigma \in \Sigma$, we denote the set $\{p \mid (q, \sigma, p) \in \delta\}$ by

$\delta(q, \sigma)$, and we can extend this notation to $w \in \Sigma^*$, and to $R \subseteq Q$. The language accepted by \mathcal{A} is $\mathcal{L}(\mathcal{A}) = \{w \in \Sigma^* \mid \delta(q_0, w) \cap F \neq \emptyset\}$.

2.1 The Glushkov Automaton

The Glushkov, or position, automaton was independently introduced by Glushkov [Glu61] and McNaughton and Yamada [MY60]. The states in the Glushkov automaton, representing a regular expression α , correspond to the positions of letters in α plus an additional initial state. Let $\bar{\alpha}$ denote the regular expression obtained by marking each letter with its position in α . The marked version of the previous example is $\bar{\tau} = (a_1 + b_2)(a_3^* + b_4a_5^* + b_6^*)^*$. The same notation is used to remove the markings, i.e., $\bar{\bar{\alpha}} = \alpha$. Now, let $\text{Pos}(\alpha) = \{1, 2, \dots, |\alpha|_\Sigma\}$ be the set of positions for $\alpha \in \mathcal{R}$, and let $\text{Pos}_0(\alpha) = \text{Pos}(\alpha) \cup \{0\}$. Then, the construction of the Glushkov automaton is based on the position sets $\text{First}(\bar{\alpha})$, $\text{Last}(\bar{\alpha})$, and $\text{Follow}(\bar{\alpha})$. These sets can be inductively defined as follows:

$$\begin{aligned} \text{First}(\emptyset) &= \text{First}(\varepsilon) = \emptyset & \text{First}(\alpha + \beta) &= \text{First}(\alpha) \cup \text{First}(\beta) \\ \text{First}(\sigma_i) &= \{i\} & \text{First}(\alpha\beta) &= \begin{cases} \text{First}(\alpha) \cup \text{First}(\beta) & \text{if } \varepsilon(\alpha) = \varepsilon \\ \text{First}(\alpha) & \text{otherwise.} \end{cases} \\ \text{First}(\alpha^*) &= \text{First}(\alpha) \end{aligned}$$

The definition of Last is almost identical and differs only for the case of concatenation, which is

$$\text{Last}(\alpha\beta) = \begin{cases} \text{Last}(\alpha) \cup \text{Last}(\beta) & \text{if } \varepsilon(\beta) = \varepsilon \\ \text{Last}(\beta) & \text{otherwise.} \end{cases}$$

The set Follow can be computed as by

$$\begin{aligned} \text{Follow}(\emptyset) &= \text{Follow}(\varepsilon) = \text{Follow}(\sigma_j) = \emptyset \\ \text{Follow}(\alpha + \beta) &= \text{Follow}(\alpha) \cup \text{Follow}(\beta) \\ \text{Follow}(\alpha\beta) &= \text{Follow}(\alpha) \cup \text{Follow}(\beta) \cup \text{Last}(\alpha) \times \text{First}(\beta) \\ \text{Follow}(\alpha^*) &= \text{Follow}(\alpha) \cup \text{Last}(\alpha) \times \text{First}(\alpha). \end{aligned}$$

The *Glushkov automaton* for α is $\mathcal{A}_{\text{pos}}(\alpha) = (\text{Pos}_0(\alpha), \Sigma, \delta_{\text{pos}}, 0, F)$, with $\delta_{\text{pos}} = \{(0, \bar{\sigma}_j, j) \mid j \in \text{First}(\bar{\alpha})\} \cup \{(i, \bar{\sigma}_j, j) \mid (i, j) \in \text{Follow}(\bar{\alpha})\}$ and $F = \text{Last}(\bar{\alpha}) \cup \{0\}$ if $\varepsilon(\alpha) = \varepsilon$, and $F = \text{Last}(\bar{\alpha})$, otherwise. Note that the number of states of $\mathcal{A}_{\text{pos}}(\alpha)$ is exactly $n+1$, where $n = |\alpha|_\Sigma$. On the other hand, the number of transitions in $\mathcal{A}_{\text{pos}}(\alpha)$ is in the worst case $n^2 + n$. Consequently, the time-complexity of any construction algorithm for $\mathcal{A}_{\text{pos}}(\alpha)$ must be at least $O(n^2)$. Considering the simplicity of the recursive definitions of the position sets used for the construction of $\mathcal{A}_{\text{pos}}(\alpha)$, an algorithm of this complexity should not be hard to find. Nevertheless, a naive implementation leads to a $O(n^3)$ algorithm, such as the one proposed by Berry and Sethi [BS86]. This is due to possibly non-disjoint unions of sets in the rule for α^* in the recursive definition of $\text{Follow}(\alpha)$. To overcome this problem, several techniques for the construction of $\mathcal{A}_{\text{pos}}(\alpha)$ were proposed over the years. The first one, of order $O(m + n^2)$, where $m = |\alpha|$,

was proposed by Brüggemann-Klein in 1993 [BK93] and it is primarily based on the prior transformation of α into *star-normal form*. Other quadratic, however sophisticated, algorithms have been introduced in 1996 and 1997, respectively in [PZC97] and [CP97].

Our goal in the next section, is to present an alternative recursive definition of $\text{Follow}(\alpha)$, that only involves disjoint unions of sets, allowing for simple implementations of that construction in time $O(n^2)$. This definition also allows us to define a cost generating function of the exact number of transitions in the Glushkov automaton in Section 4.

3 A New Algorithm for Computing $\text{Follow}(\alpha)$

In this section we define a new function E , such that for every marked regular expression α we have $\text{Follow}(\alpha) = E(\alpha)$. This function has the advantage that all unions in its definition are clearly disjoint. Our definition of E was inspired by the construction of $\mathcal{A}_{\text{pos}}(\alpha)$ by Leiss [Lei80] and shows some similarities to the transformation algorithm of α into star-normal-form by Brüggemann-Klein. Let E be given by

$$\begin{aligned} E(\emptyset) &= E(\varepsilon) = E(\sigma_i) = \emptyset \\ E(\alpha + \beta) &= E(\alpha) \cup E(\beta) \\ E(\alpha\beta) &= E(\alpha) \cup E(\beta) \cup \text{Last}(\alpha) \times \text{First}(\beta) \\ E(\alpha^*) &= E^*(\alpha) \end{aligned} \tag{2}$$

$$\begin{aligned} E^*(\emptyset) &= E^*(\varepsilon) = \emptyset \\ E^*(\sigma_i) &= \{(i, i)\} \\ E^*(\alpha + \beta) &= E^*(\alpha) \cup E^*(\beta) \cup \text{Cross}(\alpha, \beta) \\ E^*(\alpha\beta) &= \begin{cases} E^*(\alpha) \cup E^*(\beta) \cup \text{Cross}(\alpha, \beta) & \text{if } \varepsilon(\alpha) = \varepsilon(\beta) = \varepsilon \\ E^*(\alpha) \cup E(\beta) \cup \text{Cross}(\alpha, \beta) & \text{if } \varepsilon(\beta) = \varepsilon \\ E(\alpha) \cup E^*(\beta) \cup \text{Cross}(\alpha, \beta) & \text{if } \varepsilon(\alpha) = \varepsilon \\ E(\alpha) \cup E(\beta) \cup \text{Cross}(\alpha, \beta) & \text{otherwise} \end{cases} \\ E^*(\alpha^*) &= E^*(\alpha), \end{aligned} \tag{3}$$

with $\text{Cross}(\alpha, \beta) = \text{Last}(\alpha) \times \text{First}(\beta) \cup \text{Last}(\beta) \times \text{First}(\alpha)$.

Proposition 1. *For every regular expression γ we have $\text{Follow}(\bar{\gamma}) = E(\bar{\gamma})$.*

Proof.

The proof follows by induction on the structure of γ . The result is trivially true for $\bar{\gamma} = \emptyset, \varepsilon, \sigma_i, \alpha + \beta, \alpha\beta$. For $\bar{\gamma} = \delta^*$, it is sufficient to show that one has $\text{Follow}(\delta) \cup \text{Last}(\delta) \times \text{First}(\delta) = E^*(\delta)$.

For $\delta = \emptyset$, $\delta = \varepsilon$ and $\delta = \sigma_i$ this equation evaluates to $\emptyset = \emptyset$, $\emptyset = \emptyset$ and $\{(i, i)\} = \{(i, i)\}$, respectively. For $\delta = \alpha + \beta$ we have

$$\begin{aligned} & \text{Follow}(\alpha + \beta) \cup \text{Last}(\alpha + \beta) \times \text{First}(\alpha + \beta) = \\ & = (\text{Follow}(\alpha) \cup \text{Last}(\alpha) \times \text{First}(\alpha)) \cup (\text{Follow}(\beta) \cup \text{Last}(\beta) \times \text{First}(\beta)) \cup \\ & \quad \cup \text{Last}(\alpha) \times \text{First}(\beta) \cup \text{Last}(\beta) \times \text{First}(\alpha) = \\ & = \mathbf{E}^*(\alpha) \cup \mathbf{E}^*(\beta) \cup \text{Last}(\alpha) \times \text{First}(\beta) \cup \text{Last}(\beta) \times \text{First}(\alpha) = \mathbf{E}^*(\alpha + \beta). \end{aligned}$$

We illustrate the proof for $\delta = \alpha\beta$ with the case where $\varepsilon(\alpha) \neq \varepsilon$ and $\varepsilon(\beta) = \varepsilon$:

$$\begin{aligned} & \text{Follow}(\alpha\beta) \cup \text{Last}(\alpha\beta) \times \text{First}(\alpha\beta) = \\ & = \text{Follow}(\alpha) \cup \text{Follow}(\beta) \cup \text{Last}(\alpha) \times \text{First}(\beta) \cup \\ & \quad \cup \text{Last}(\alpha) \times \text{First}(\alpha) \cup \text{Last}(\beta) \times \text{First}(\alpha) \\ & = \mathbf{E}^*(\alpha) \cup \mathbf{E}(\beta) \cup \text{Last}(\alpha) \times \text{First}(\beta) \cup \text{Last}(\beta) \times \text{First}(\alpha) = \mathbf{E}^*(\alpha\beta). \end{aligned}$$

Finally, for $\delta = \alpha^*$ we have

$$\begin{aligned} & \text{Follow}(\alpha^*) \cup \text{Last}(\alpha^*) \times \text{First}(\alpha^*) = \\ & = \text{Follow}(\alpha) \cup \text{Last}(\alpha) \times \text{First}(\alpha) \cup \text{Last}(\alpha) \times \text{First}(\alpha) \\ & = \text{Follow}(\alpha) \cup \text{Last}(\alpha) \times \text{First}(\alpha) = \mathbf{E}^*(\alpha) = \mathbf{E}^*(\alpha^*). \quad \square \end{aligned}$$

4 Counting the Number of Transitions in the Glushkov Automaton

Nicaud [Nic09] showed that the average number of transitions in the Glushkov automaton $\mathcal{A}_{\text{pos}}(\alpha)$ is $O(|\alpha|)$. However, his computation of the number of transitions was not exact because the definition used for the Follow function did not take into account the possible non-disjoint unions of its results. In this section, based on the algorithm E we compute the exact number of transitions in $\mathcal{A}_{\text{pos}}(\alpha)$, $E_k(z)$, as well as its average cardinality, $T_k(z)$. This is done by the use of the standard methods of analytic combinatorics as expounded by Flajolet and Sedgewick [FS08]. These apply to generating functions $A(z) = \sum_n a_n z^n$ for a combinatorial class \mathcal{A} with a_n objects of size n , or cost generating functions $C(z) = \sum_\alpha c(\alpha) z^{|\alpha|}$, where $c(\alpha)$ is some measure of the object $\alpha \in \mathcal{A}$.

In this section we compute and study the cost generating functions $E_k(z)$ and $T_k(z)$, and their asymptotic behaviours. The other functions used herein, as well as details on how to obtain them, can be found in the above cited article and in Broda *et al* [BMMR11b]. A more detailed description of the below computations can be found in a companion technical report of this paper [BMMR11a].

4.1 The Average Number of Transitions in $\mathcal{A}_{\text{pos}}(\alpha)$

For counting purposes, we will consider regular expressions as defined in (1), but without \emptyset . Note that this limitation only excludes the empty language.

The functions that count the cardinalities of $\text{First}(\bar{\alpha})$, $\text{Last}(\bar{\alpha})$, and $\text{E}(\bar{\alpha})$, are respectively denoted by $f(\alpha)$, $l(\alpha)$, and $e(\alpha)$. Given the definitions of $f(\alpha)$ and $l(\alpha)$, $e(\alpha)$ satisfies the following:

$$\begin{aligned} e(\sigma) &= e(\varepsilon) = 0, \\ e(\alpha + \beta) &= e(\alpha) + e(\beta), \\ e(\alpha\beta) &= e(\alpha) + e(\beta) + l(\alpha) \cdot f(\beta), \\ e(\alpha^*) &= e^*(\alpha), \end{aligned} \tag{4}$$

where $e^*(\alpha)$ is given by,

$$\begin{aligned} e^*(\varepsilon) &= 0, \quad e^*(\sigma) = 1, \\ e^*(\alpha + \beta) &= e^*(\alpha) + e^*(\beta) + c(\alpha, \beta), \\ e^*(\alpha_\varepsilon\beta_\varepsilon) &= e^*(\alpha_\varepsilon) + e^*(\beta_\varepsilon) + c(\alpha_\varepsilon, \beta_\varepsilon), \\ e^*(\alpha_{\bar{\varepsilon}}\beta_{\bar{\varepsilon}}) &= e^*(\alpha_{\bar{\varepsilon}}) + e^*(\beta_{\bar{\varepsilon}}) + c(\alpha_{\bar{\varepsilon}}, \beta_{\bar{\varepsilon}}), \\ e^*(\alpha_\varepsilon\beta_{\bar{\varepsilon}}) &= e(\alpha_\varepsilon) + e^*(\beta_{\bar{\varepsilon}}) + c(\alpha_\varepsilon, \beta_{\bar{\varepsilon}}), \\ e^*(\alpha_{\bar{\varepsilon}}\beta_\varepsilon) &= e(\alpha_{\bar{\varepsilon}}) + e(\beta_\varepsilon) + c(\alpha_{\bar{\varepsilon}}, \beta_\varepsilon), \\ e^*(\alpha^*) &= e^*(\alpha). \end{aligned} \tag{5}$$

with $c(\alpha, \beta) = l(\alpha) \cdot f(\beta) + l(\beta) \cdot f(\alpha)$. Then, the function

$$t(\alpha) = f(\alpha) + e(\alpha)$$

computes the number of transitions in the Glushkov automaton of α . The cost generating function associated to t is given by

$$T_k(z) = F_k(z) + E_k(z),$$

where $F_k(z)$ and $E_k(z)$ are the cost generating functions associated to f and e , respectively. By symmetry, the cost generating function $L_k(z)$ associated to l is the same as $F_k(z)$, *i.e.*

$$L_k(z) = F_k(z) = \frac{kz}{1 - z - 3zR_k(z) - zR_{k,\varepsilon}(z)}.$$

In this last expression $R_k(z)$ and $R_{k,\varepsilon}(z)$ denote respectively the generating functions for regular expressions, and for regular expressions whose languages contain ε and are given by

$$R_k(z) = \frac{1 - z - \sqrt{\Delta_k(z)}}{4z} \quad \text{and} \quad R_{k,\varepsilon}(z) = \frac{z + zR_k(z)}{1 - 2zR_k(z)}, \tag{6}$$

where $\Delta_k(z) = 1 - 2z - (7 + 8k)z^2$. Hence, for the number of regular expressions of size n , one has

$$[z^n]R_k(z) \sim \frac{\sqrt{2(1 - \rho_k)}}{8\rho_k\sqrt{\pi}} \rho_k^{-n} n^{-3/2}, \quad \text{where } \rho_k = \frac{1}{1 + \sqrt{8k + 8}}. \tag{7}$$

From the equations in (4) one can compute the associated cost generating functions $E_k(z)$ and $E_k^*(z)$. For instance, the equation for concatenation contributes with the term $2zE_k(z)R_k(z) + zF_k(z)^2$ in the equation for $E_k(z)$. Collecting all terms the following equations must be satisfied

$$E_k(z) = 4zE_k(z)R_k(z) + zF_k(z)^2 + zE_k^*(z)$$

$$E_k^*(z) = kz + 2zE_k^*(z)R_k(z) + 2zE_k^*(z)R_{k,\varepsilon}(z) + 4zF_k(z)^2 + 2zE_k(z)R_{k,\bar{\varepsilon}}(z) + zE_k^*(z).$$

After simplification one gets

$$E_k(z) = \frac{kz^2 + zF_k(z)^2\Lambda_k(z) + 4z^2F_k(z)^2}{(1 - 4zR_k(z))\Lambda_k(z) - 2z^2R_{k,\bar{\varepsilon}}(z)}, \quad (8)$$

where $\Lambda_k(z) = 1 - z - 2zR_{k,\varepsilon}(z) - 2zR_k(z)$. After substituting the functions in (8) by their expressions in terms of z and k , one obtains

$$T_k(z) = \frac{P_k(z)}{Q_k(z)\sqrt{\Delta_k(z)}}, \quad (9)$$

where $P_k(z)$ is a polynomial in $\sqrt{\Delta_k(z)}$ over $\mathbb{Q}[z]$, with $Q_k(z)$ given by

$$Q_k(z) = \left(1 - 2z - 7z^2 + 4(1+z)\sqrt{\Delta_k(z)} + 3\Delta_k(z)\right)^2 \\ \left(1 - 5z^2 + 2(1+2z)\sqrt{\Delta_k(z)} + \Delta_k(z)\right).$$

This function $Q_k(z)$ is positive for all values of z in the real segment $[0, \rho_k]$, because $1 - 2z - 7z^2 = 8kz^2 + \Delta_k(z)$ and $1 - 5z^2 = 2z + 2z^2 + 8kz^2 + \Delta_k(z)$, and $\Delta_k(z)$ is non-negative in that segment. By Pringsheim's Theorem (Theorem IV.6 of [FS08], p. 240) one can conclude that $T_k(z)$ has radius of convergence equal to ρ_k . Moreover, it can be shown that $T_k(z)$ has no singularities on the the boundary of its disc of convergence, $\|z\| = \rho_k$, besides the one at $z = \rho_k$.

Using exactly the same technique employed in the previously referred article, one obtains

$$T_k(z) = \frac{P_k(\rho_k)}{\sqrt{2-2\rho_k} Q_k(\rho_k)} \frac{1}{\sqrt{1-z/\rho_k}} + o\left(\frac{1}{\sqrt{1-z/\rho_k}}\right), \quad (10)$$

from which it follows that

$$[z^n]T_k(z) \sim \frac{P_k(\rho_k)}{\sqrt{\pi} \sqrt{2-2\rho_k} Q_k(\rho_k)} \rho_k^{-n} n^{-\frac{1}{2}}. \quad (11)$$

Using the actual expression of P_k , which we omit due to lack of space, one can get

$$[z^n]T_k(z) \sim \frac{(1+\rho_k)(2+16\rho_k+10\rho_k^2-12\rho_k^3)}{8\rho_k\sqrt{\pi}(1-5\rho_k^2)\sqrt{2-2\rho_k}} \rho_k^{-n} n^{-\frac{1}{2}}. \quad (12)$$

Considering the cost generating function for the number of letters in an element $\alpha \in \mathcal{R}$, computed by Nicaud to be equal to

$$Let_k(z) = \frac{kz}{\sqrt{\Delta_k(z)}},$$

and for which

$$[z^n]Let_k(z) \sim \frac{k\rho_k}{\sqrt{\pi(2-2\rho_k)}} \rho_k^{-n} n^{-1/2},$$

one gets an asymptotic expression for the average number of transitions per state:

$$\frac{[z^n]T_k(z)}{[z^n]Let_k(z)} \sim \frac{(1+\rho_k)(2+16\rho_k+10\rho_k^2-12\rho_k^3)}{(1-2\rho_k-7\rho_k^2)(1-5\rho_k^2)}. \quad (13)$$

And finally, one has for the average number of transitions per regular expression the following asymptotic estimation:

$$\frac{[z^n]T_k(z)}{[z^n]R_k(z)} \sim \frac{(1+\rho_k)(1+8\rho_k+5\rho_k^2-6\rho_k^3)}{(1-\rho_k)(1-5\rho_k^2)} n. \quad (14)$$

Since ρ_k tends to 0 as k goes to ∞ , it follows that for large k the average number of transitions per state is approximately 2, while the average number of transitions per automaton is approximately the size of the original regular expression.

5 The Average Number of Transitions in \mathcal{A}_{pd}

The partial derivative automaton $\mathcal{A}_{pd}(\alpha)$ of a regular expression α was defined independently by Mirkin's [Mir66] and Antimirov [Ant96]. Champarnaud and Ziadi stated the equivalence of the two formulations [CZ01], and proved that \mathcal{A}_{pd} is a quotient of the Glushkov automaton \mathcal{A}_{pos} [CZ02]. This means that $\mathcal{A}_{pd}(\alpha)$ can be obtained from $\mathcal{A}_{pos}(\alpha)$ by the merging of states belonging to the same equivalence class. That, on the other hand, may lead to the merging of transitions. In this section, we estimate the average number of transitions of $\mathcal{A}_{pd}(\alpha)$ when compared with the ones of $\mathcal{A}_{pos}(\alpha)$. For this, it is essential to have the exact counting of the number of transitions of $\mathcal{A}_{pos}(\alpha)$ obtained in Section 4.1.

The $\mathcal{A}_{pd}(\alpha)$ can be defined using the notion of partial derivative, introduced by Antimirov as a non-deterministic version of Brzozowski's derivative [Brz64].

For a regular expression α and a letter $\sigma \in \Sigma$, the set $\partial_\sigma(\alpha)$ of *partial derivatives* of α w.r.t. σ is defined inductively as follows:

$$\begin{aligned} \partial_\sigma(\emptyset) &= \partial_\sigma(\varepsilon) = \emptyset & \partial_\sigma(\alpha + \beta) &= \partial_\sigma(\alpha) \cup \partial_\sigma(\beta) \\ \partial_\sigma(\sigma') &= \begin{cases} \{\varepsilon\}, & \text{if } \sigma' = \sigma \\ \emptyset, & \text{otherwise} \end{cases} & \partial_\sigma(\alpha_\varepsilon\beta) &= \partial_\sigma(\alpha_\varepsilon)\beta \cup \partial_\sigma(\beta) \\ \partial_\sigma(\alpha^*) &= \partial_\sigma(\alpha)\alpha^* & \partial_\sigma(\alpha_{\bar{\varepsilon}}\beta) &= \partial_\sigma(\alpha_{\bar{\varepsilon}})\beta \end{aligned}$$

where for any $S \subseteq \mathcal{R}$, $S\emptyset = \emptyset S = \emptyset$, and $S\varepsilon = \varepsilon S = S$. This definition can be extended to sets of regular expressions, to words, and to languages in the obvious way. The *set of partial derivatives* of α , $\{\partial_w(\alpha) \mid w \in \Sigma^*\}$, is denoted by $P(\alpha)$. The *partial derivative automaton* $\mathcal{A}_{\text{pd}}(\alpha)$ is defined by $\mathcal{A}_{\text{pd}}(\alpha) = (P(\alpha), \Sigma, \delta_{\text{pd}}, \alpha, \{q \in P(\alpha) \mid \varepsilon(q) = \varepsilon\})$, where $\delta_{\text{pd}}(q, \sigma) = \partial_\sigma(q)$, for all $q \in P(\alpha)$ and $\sigma \in \Sigma$. Antimirov proved that $\mathcal{L}(\mathcal{A}_{\text{pd}}(\alpha)) = \mathcal{L}(\alpha)$.

Using Mirkin's formulation one has $P(\alpha) = \pi(\alpha) \cup \{\alpha\}$, where the set $\pi(\alpha)$ is inductively defined as follows:

$$\begin{aligned} \pi(\emptyset) &= \emptyset & \pi(\alpha + \beta) &= \pi(\alpha) \cup \pi(\beta) \\ \pi(\varepsilon) &= \emptyset & \pi(\alpha\beta) &= \pi(\alpha)\beta \cup \pi(\beta) \\ \pi(\sigma) &= \{\varepsilon\} & \pi(\alpha^*) &= \pi(\alpha)\alpha^*. \end{aligned} \tag{15}$$

5.1 Counting the Mergings of Transitions

Broda *et al.* [BMMR11b] gave a lower bound of the number of mergings of states in $\pi(\alpha)$ with respect to $\text{Pos}(\alpha)$, which allowed to obtain an upper bound on the average state complexity of $\mathcal{A}_{\text{pd}}(\alpha)$. There, it was observed that the merging of states is primarily caused by sub-expressions γ of α such that $\varepsilon \in \pi(\gamma)$. In fact, in the presence of sub-expressions with this property, denoted by α_{π_ε} , during the computation of $\pi(\alpha)$ some unions may not be disjoint.

In this section, and using the same technique, we determine a lower bound of the number of mergings of transitions. Considering (2) and in particular, the concatenation case, it is easy to see, that whenever there is a merging of two states in the set $\text{Last}(\alpha)$, there are exactly $f(\beta) = |\text{First}(\beta)|$ mergings of transitions. Although there can be merging of states of $\text{First}(\beta)$, they will not be considered in the computation of that lower bound. We first compute a lower bound for the number of mergings $i_\ell(\alpha)$ of states i such that $i \in \text{Last}(\alpha)$.

In addition to α_ε and α_{π_ε} , we use the subclass of regular expressions $\alpha_{r,\varepsilon}$ such that $\alpha_{r,\varepsilon} \in \pi(\alpha_{r,\varepsilon})$ and $\varepsilon(\alpha_{r,\varepsilon}) = \varepsilon$.

The grammar for α_{π_ε} and its generating function $R_{k,\pi_\varepsilon}(z)$ are the ones used by Broda *et al.* For $\alpha_{r,\varepsilon}$ one has

$$\alpha_{r,\varepsilon} := \alpha_{\pi_\varepsilon}^* \mid \alpha_{r,\varepsilon} \cdot \alpha_\varepsilon, \quad \text{and} \quad R_{k,r,\varepsilon}(z) = \frac{zR_{k,\pi_\varepsilon}(z)}{1 - zR_{k,\varepsilon}(z)}.$$

Finally, in order to get a better approximation for the counting of the state mergings, we define $i_\ell(\alpha)$ by the following:

$$\begin{aligned} i_\ell(\emptyset) &= i_\ell(\varepsilon) = i_\ell(\sigma) = 0, & i_\ell(\alpha_{\pi_\varepsilon} \alpha_{r,\varepsilon}) &= i_\ell(\alpha_{\pi_\varepsilon}) + i_\ell(\alpha_{r,\varepsilon}) + 1, \\ i_\ell(\alpha_{\pi_\varepsilon} + \alpha_{\pi_\varepsilon}) &= i_\ell(\alpha_{\pi_\varepsilon}) + i_\ell(\alpha_{\pi_\varepsilon}) + 1, & i_\ell(\alpha_{\pi_\varepsilon} \alpha_{r,\bar{\varepsilon}}) &= i_\ell(\alpha_{r,\bar{\varepsilon}}), \\ i_\ell(\alpha_{\pi_\varepsilon} + \alpha_{\bar{\pi}_\varepsilon}) &= i_\ell(\alpha_{\pi_\varepsilon}) + i_\ell(\alpha_{\bar{\pi}_\varepsilon}), & i_\ell(\alpha_{\pi_\varepsilon} \alpha_{\bar{r},\varepsilon}) &= i_\ell(\alpha_{\pi_\varepsilon}) + i_\ell(\alpha_{\bar{r},\varepsilon}), \\ i_\ell(\alpha_{\bar{\pi}_\varepsilon} + \alpha) &= i_\ell(\alpha_{\bar{\pi}_\varepsilon}) + i_\ell(\alpha), & i_\ell(\alpha_{\pi_\varepsilon} \alpha_{\bar{r},\bar{\varepsilon}}) &= i_\ell(\alpha_{\bar{r},\bar{\varepsilon}}), \\ i_\ell(\alpha^*) &= i_\ell(\alpha), & i_\ell(\alpha_{\bar{\pi}_\varepsilon} \alpha_\varepsilon) &= i_\ell(\alpha_{\bar{\pi}_\varepsilon}) + i_\ell(\alpha_\varepsilon), \\ & & i_\ell(\alpha_{\bar{\pi}_\varepsilon} \alpha_{\bar{\varepsilon}}) &= i_\ell(\alpha_{\bar{\varepsilon}}), \end{aligned}$$

where $\gamma_{\bar{x}}$ denotes the complement of γ_x .

To illustrate the previous rules, consider the case of $i_\ell(\alpha_{\pi_\varepsilon}\alpha_{r,\varepsilon})$. Here, one has $\pi(\alpha_{\pi_\varepsilon}\alpha_{r,\varepsilon}) = \pi(\alpha_{\pi_\varepsilon})\alpha_{r,\varepsilon} \cup \pi(\alpha_{r,\varepsilon})$. By definition, $\varepsilon \in \pi(\alpha_{\pi_\varepsilon})$ and $\alpha_{r,\varepsilon} \in \pi(\alpha_{r,\varepsilon})$. Hence $\alpha_{r,\varepsilon}$ belongs to both $\pi(\alpha_{\pi_\varepsilon})\alpha_{r,\varepsilon}$ and $\pi(\alpha_{r,\varepsilon})$, which causes a merging of two states in $\text{Last}(\alpha)$. This merging is accounted for by the 1 in the definition of $i_\ell(\alpha_{\pi_\varepsilon}\alpha_{r,\varepsilon})$. On the other hand, since $\varepsilon(\alpha_{r,\varepsilon}) = \varepsilon$, one also has to count all mergings of states in $\text{Last}(\alpha_{\pi_\varepsilon})$ (counted by $i_\ell(\alpha_{\pi_\varepsilon})$), besides those in $\text{Last}(\alpha_{r,\varepsilon})$.

The generating function, $I_k(z)$, of i_ℓ satisfies the following:

$$I_k(z) = \frac{zR_{k,\pi_\varepsilon}(z)^2 + zR_{k,\pi_\varepsilon}(z)R_{k,r,\varepsilon}(z)}{1 - z - 3zR_k(z) - zR_{k,\varepsilon}(z)},$$

where $R_{k,\pi_\varepsilon}(z) = \frac{z^2 + 3zR_k(z) - 1 + \sqrt{(z^2 + 3zR_k(z) - 1)^2 + 4kz^2}}{2z}$.

Using (2) and (3), one can easily define a function that computes a lower bound for the number of transition mergings in the Glushkov automaton:

$$\begin{aligned} i_t(\varepsilon) &= i_t(\sigma) = 0 & i_t^*(\alpha + \beta) &= i_t^*(\alpha) + i_t^*(\beta) + \mathbf{c}_t(\alpha, \beta) \\ i_t(\alpha + \beta) &= i_t(\alpha) + i_t(\beta) & i_t^*(\alpha_\varepsilon\beta_\varepsilon) &= i_t^*(\alpha_\varepsilon) + i_t^*(\beta_\varepsilon) + \mathbf{c}_t(\alpha_\varepsilon, \beta_\varepsilon) \\ i_t(\alpha\beta) &= i_t(\alpha) + i_t(\beta) + i_\ell(\alpha) \cdot \mathbf{f}(\beta) & i_t^*(\alpha_{\bar{\varepsilon}}\beta_\varepsilon) &= i_t^*(\alpha_{\bar{\varepsilon}}) + i_t(\beta_\varepsilon) + \mathbf{c}_t(\alpha_{\bar{\varepsilon}}, \beta_\varepsilon) \\ i_t(\alpha^*) &= i_t^*(\alpha) & i_t^*(\alpha_\varepsilon\beta_{\bar{\varepsilon}}) &= i_t(\alpha_\varepsilon) + i_t^*(\beta_{\bar{\varepsilon}}) + \mathbf{c}_t(\alpha_\varepsilon, \beta_{\bar{\varepsilon}}) \\ i_t^*(\varepsilon) &= 0 & i_t^*(\alpha_{\bar{\varepsilon}}\beta_{\bar{\varepsilon}}) &= i_t(\alpha_{\bar{\varepsilon}}) + i_t(\beta_{\bar{\varepsilon}}) + \mathbf{c}_t(\alpha_{\bar{\varepsilon}}, \beta_{\bar{\varepsilon}}) \\ i_t^*(\sigma) &= 1 & i_t^*(\alpha^*) &= i_t^*(\alpha) \end{aligned}$$

where $\mathbf{c}_t(\alpha, \beta) = i_\ell(\alpha) \cdot \mathbf{f}(\beta) + i_\ell(\beta) \cdot \mathbf{f}(\alpha)$. The corresponding generating function satisfies the following:

$$It_k(z) = \frac{z\Lambda_k(z)I_k(z)F_k(z) + kz^2 + 4z^2I_k(z)F_k(z)}{(1 - 4zR_k(z))\Lambda_k(z) - 2z^2R_{k,\bar{\varepsilon}}(z)},$$

where $\Lambda_k(z) = 1 - z - 2zR_k(z) - 2zR_{k,\varepsilon}(z)$. Analogously to what was done before, one has

$$[z^n]It_k(z) \sim \frac{(1 + \rho_k)(a(\rho_k)b(\rho_k) + c(\rho_k))}{16\sqrt{\pi}\rho_k\sqrt{2 - 2\rho_k}(1 - 5\rho_k^2)d(\rho_k)}\rho_k^{-n}n^{-\frac{1}{2}} \quad (16)$$

where

$$\begin{aligned} a(z) &= -2 - 23z - 77z^2 - 50z^3 + 92z^4 + 77z^5 - 13z^6 - 4z^7 \\ b(z) &= \sqrt{9 - 10z - 55z^2 - 24z^3 + 16z^4} \\ c(z) &= 10 + 89z + 54z^2 - 603z^3 - 1114z^4 - 349z^5 + \\ &\quad + 130z^6 - 209z^7 - 40z^8 - 16z^9 \\ d(z) &= (1 - 2z - 7z^2)(2 + z - 3z^2). \end{aligned}$$

Therefore, a lower bound for the average number of mergings per transition of the Glushkov automaton is given by

$$[z^n]\frac{It_k(z)}{T_k(z)} \sim \frac{a(\rho_k)b(\rho_k) + c(\rho_k)}{4(1 + 8\rho_k + 5\rho_k^2 - 6\rho_k^3)d(\rho_k)}. \quad (17)$$

This means that, asymptotically with respect to k , the number of transitions in \mathcal{A}_{pd} is at most half the number of transitions in \mathcal{A}_{pos} .

6 Comparison with Experimental Results

We compared the estimates obtained in the previous sections with some experimental results. For each $k \in \{2, 3, 10, 30, 50\}$, the experiment consisted of the comparison of the sizes of \mathcal{A}_{pos} and \mathcal{A}_{pd} , that were computed for each regular expression in the samples of 1000 uniform random generated regular expressions of size 1000. Table 1 presents the average values obtained, and columns eight and ten give the asymptotic ratios obtained in (13) and in (17), respectively. The quality of the approximation of the asymptotic average number of transitions per state for \mathcal{A}_{pos} , and that the actual values are close to the limit even for relatively small alphabets is evident from the table. The upper bound for the ratio of the transition complexity of \mathcal{A}_{pd} to the one of \mathcal{A}_{pos} is within an error less than 15%. The experimental values also suggest that the number of transitions of \mathcal{A}_{pd} is on average, and as k grows, the alphabetic size of the original regular expression.

k	$ \alpha $	$ \alpha _{\Sigma}$	$ \delta_{pos} $	$ P(\alpha) $	$ \delta_{pd} $	$\frac{ \delta_{pos} }{ \alpha _{\Sigma}+1}$	$\frac{[z^n]T_k(z)}{[z^n]Let_k(z)+1}$	$\frac{ \delta_{pd} }{ \delta_{pos} }$	$1 - \frac{[z^n]It_k(z)}{[z^n]T_k(z)}$
2	1000	276	3345	187	1806	12.1	12.2	0.54	0.68
3	1000	318	2997	206	1564	9.4	9.6	0.52	0.64
10	1000	405	2203	236	1079	5.4	5.3	0.49	0.58
30	1000	453	1676	247	796	3.7	3.6	0.47	0.54
50	1000	466	1516	250	718	3.3	3.2	0.47	0.53
100	—	—	—	—	—	—	2.8	—	0.53
1000	—	—	—	—	—	—	2.2	—	0.49

Table 1. Experimental results for uniform random generated regular expressions

7 Conclusions

In this paper we presented a new algorithm for computing Follow, which is quadratic in the alphabetic size of the original regular expression, and that leads to a straightforward and direct implementation. This algorithm allowed us to exactly count the number of transitions of the Glushkov automaton. Using this, we computed the average number of transitions of that automaton, concluding that, for large alphabets, it is approximately the double of the original regular expression alphabetic size.

Considering special sub-classes of regular expressions that are primarily responsible for state mergings, we computed an upper bound for the number of transitions in the partial derivative automaton. We, then, used analytic combinatorial methods to obtain average values and asymptotic limits for that number,

concluding that, on average and asymptotically, the partial derivative automaton has at most half the number of transitions of the Glushkov's. Experimental figures corroborate these results.

References

- [Ant96] V. M. Antimirov. Partial derivatives of regular expressions and finite automaton constructions. *Theoret. Comput. Sci.*, 155(2):291–319, 1996.
- [BK93] A. Brüggemann-Klein. Regular expressions into finite automata. *Theoret. Comput. Sci.*, 48:197–213, 1993.
- [BMMR11a] S. Broda, A. Machiavelo, N. Moreira, and R. Reis. On the average size of Glushkov and partial derivative automata. Technical Report DCC-2011-03, FCUP & CMUP, Universidade do Porto, April 2011.
- [BMMR11b] S. Broda, A. Machiavelo, N. Moreira, and R. Reis. On the average state complexity of partial derivative automata. *International Journal of Foundations of Computer Science*, 2011. Accepted to publication.
- [Brz64] J. A. Brzozowski. Derivatives of regular expressions. *JACM*, 11(4):481–494, October 1964.
- [BS86] G. Berry and R. Sethi. From regular expressions to deterministic automata. *Theoret. Comput. Sci.*, 48:117–126, 1986.
- [CP97] C.-H. Chang and R. Paige. From regular expressions to DFA's using compressed NFA's. *Theor. Comput. Sci.*, 178(1-2):1–36, 1997.
- [CZ01] J. M. Champarnaud and D. Ziadi. From Mirkin's prebases to Antimirov's word partial derivatives. *Fundam. Inform.*, 45(3):195–205, 2001.
- [CZ02] J. M. Champarnaud and D. Ziadi. Canonical derivatives, partial derivatives and finite automaton constructions. *Theoret. Comput. Sci.*, 289:137–163, 2002.
- [FS08] P. Flajolet and R. Sedgewick. *Analytic Combinatorics*. Cambridge University Press, 2008.
- [Glu61] V. M. Glushkov. The abstract theory of automata. *Russian Math. Surveys*, 16:1–53, 1961.
- [HK10] M. Holzer and M. Kutrib. The complexity of regular(-like) expressions. In Y. Gao, H. Lu, S. Seki, and S. Yu, editors, *Proc. 14th DLT 2010*, volume 6224 of *LNCS*, pages 16–30. Springer, 2010.
- [Koz97] D. C. Kozen. *Automata and Computability*. Springer, 1997.
- [Lei80] E. Leiss. Constructing a finite automaton for a given regular expression. *SIGACT News*, 12(3), Sep 1980.
- [Mir66] B. G. Mirkin. An algorithm for constructing a base in a language of regular expressions. *Engineering Cybernetics*, 5:51–57, 1966.
- [MY60] R. McNaughton and H. Yamada. Regular expressions and state graphs for automata. *IEEE Transactions on Electronic Computers*, 9:39–47, 1960.
- [Nic09] C. Nicaud. On the average size of Glushkov's automata. In A. Dediu, A.-M. Ionescu, and C. M. Vide, editors, *3rd LATA 2009*, volume 5457 of *LNCS*, pages 626–637. Springer, 2009.
- [PZC97] J.-L. Ponty, D. Ziadi, and J.-M. Champarnaud. A new quadratic algorithm to convert a regular expression into an automaton. In D. R. Raymond, D. Wood, and S. Yu, editors, *Proc. 1st WIA '96*, volume 1260 of *LNCS*, pages 109–119. Springer, 1997.
- [Sak09] J. Sakarovitch. *Elements of Automata Theory*. CUP, 2009.