

# Регулярные выражения с повторными переменными. Обзор формализмов и алгоритмы сопоставления.

Выполнила: Исмаилова Д.Н., МГТУ им. Баумана

Научный руководитель: Непейвода А.Н., МГТУ им. Баумана,

ИПС им. А.К. Айламазяна РАН

Совместное Совещание по языку Рефал  
МГТУ им. Баумана (кафедра ИУ9) и  
ИПС им. А.К. Айламазяна РАН

# Проблематика

Расширенные регулярные выражения мало изучены, а потому не имеет достаточной научной базы для быстрого сопоставления с ними.

Кроме того не существует эффективных инструментов для работы с расширенными регулярными выражениями.

# План доклада

- Обзор формализмов расширенных регулярных выражений и их свойства
- Эффективный разбор академических регулярных выражений на примере Re2
- Расширение существующего механизма на расширенные регулярные выражения
- Обращение классических регулярных выражений
- Ввод формализма, удобного для преобразований регулярных выражений
- Обращение расширенных регулярных выражений
- Результаты

# Регулярные выражения с обратными ссылками

---

$(a \mid b)^*c$  - классическое регулярное выражение

$(a^*)b \backslash 1 b \backslash 1$  - регулярное выражение с обратными ссылками, соответствующее языку  $\{a^n b a^n b a^n\}$ , не являющимся даже контекстно-свободным.

---

# Алгебра Клини

**Алгебра Клини** — это полукольцо  $\langle \mathcal{A}, +, \cdot, \emptyset, 1 \rangle$ , содержащее дополнительную операцию  $*$ , идемпотентное по  $+$ , и удовлетворяющее следующим аксиомам:

- $\forall a \in \mathcal{A} (1 + a \cdot a^* = a^* \ \& \ 1 + a^* \cdot a = a^*)$  (закон раскрытия итерации слева и справа);
- $\forall a, x \in \mathcal{A} ((a \cdot x + x = x \Rightarrow a^* \cdot x + x = x) \ \& \ (x \cdot a + x = x \Rightarrow x \cdot a^* + x = x))$  (левая и правая лемма Ардена).

---

$$a(ba)^* = (ab)^*a \text{ - sliding}$$

$$a^*(ba^*)^* = (a + b)^* \text{ - denesting}$$

# Семантики

Если расширенное регулярное выражение находится в  $\epsilon$ -семантике, то неинициализированные при разборе ссылки заменяются на  $\epsilon$ .

При  $\emptyset$ -семантике слово, сопоставляющееся по пути разбора с неинициализированными ссылками, считается не соответствующим регулярному выражению.

---

$a \mid b \backslash 1$  — синтаксически некорректно

$(a \mid b \backslash 1)$  сопоставляется только с  $a$ .

# PCRE2

Скобочные группы могут быть именованными или безымянными.

- в регулярном выражении не может быть ссылки на несуществующую скобочную группу;
- скобочная группа не может входить в регулярное выражение позже, чем ссылка на её номер, за исключением групп 1–7;
- $\emptyset$ -семантика.

---

*Выражение  $(ab)c\backslash 1$  распознаёт язык  $\{abscab\}$ , тогда как выражение  $a(bc)\backslash 1$  распознаёт язык  $\{abcbcb\}$ .*

# Формализм Кампенау-Саломеа-Ю

Все скобочные группы нумеруются автоматически по первому вхождению открывающей скобки.

- Каждое появление обратной ссылки должно быть предварено соответствующей закрытой скобочной группой.
  - $\varepsilon$ -семантика.
- 

*Выражение  $(a^* \mid b^*) \setminus 1$  распознаёт язык  $\{a^{2n}\} \cup \{b^{2m}\}$ , тогда как выражение  $((a^*) \setminus 1 \mid (b^*) \setminus 1)$  некорректно в текущем формализме.*



# Формализм Шмидта

Reference word (ref-word) над алфавитом  $\Sigma$  — это выражение над алфавитом  $\Sigma \cup \{[\mathbf{x}_i, ]_{\mathbf{x}_i}, \mathbf{x}_i \mid i \in \mathbb{N}\}$ , где  $\mathbf{x}_i$  — переменная, а  $[\mathbf{x}_i, ]_{\mathbf{x}_i}$  — скобки, выделяющие подвыражение для переменной  $\mathbf{x}_i$ . При этом если выражение  $[\mathbf{x}_i \omega]_{\mathbf{x}_i}$  — ref-word, то  $\omega$  не содержит  $\mathbf{x}_i$ .

---

$$([\mathbf{x}_1 \mathbf{a} [\mathbf{x}_2 \mathbf{b}]_{\mathbf{x}_1} \mathbf{c}]_{\mathbf{x}_2})$$

# Регулярные выражения над рекурсивными образцами

Регулярные выражения над алфавитом с переменными, каждая из которых в свою очередь соответствует регулярному выражению над рекурсивными образцами.

---

$$cy^*cy^*, y = bx^*, x = a^*c(ba^n)^m(ba^n)^m$$

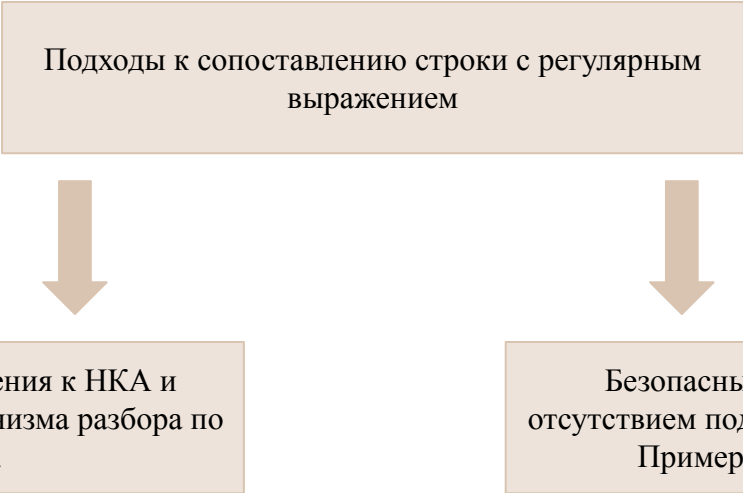
---

$$(xax)^*, x \in b^*$$

$$(x_1ax_1)^*x_2ax_2\dots x_iax_i\dots, x_i \in b^*$$

# Re2 vs обратных ссылок

Подходы к сопоставлению строки с регулярным выражением

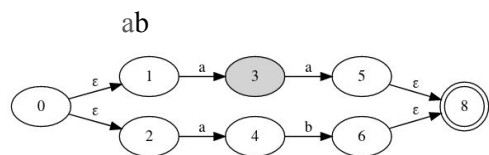
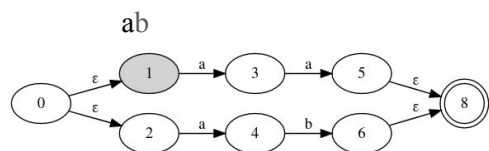


```
graph TD; A[Подходы к сопоставлению строки с регулярным выражением] --> B[Переход от регулярного выражения к НКА и последующее использование механизма разбора по НКА с возвратами.]; A --> C[Безопасная реализация без возвратов, с отсутствием поддержки расширенного синтаксиса. Примеры: модули в Go и Rust, Re2.];
```

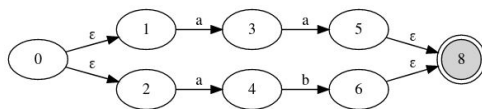
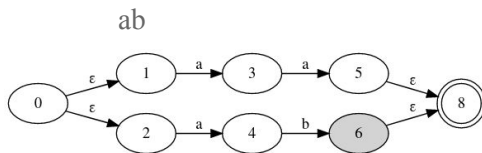
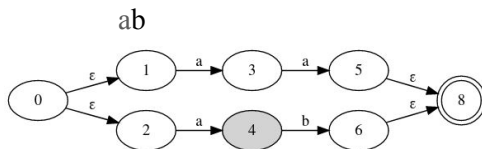
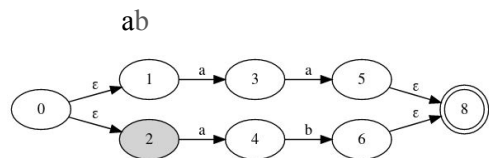
Переход от регулярного выражения к НКА и последующее использование механизма разбора по НКА с возвратами.

Безопасная реализация без возвратов, с отсутствием поддержки расширенного синтаксиса.  
Примеры: модули в *Go* и *Rust*, *Re2*.

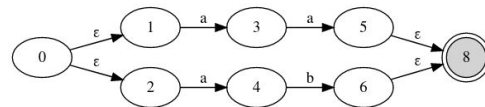
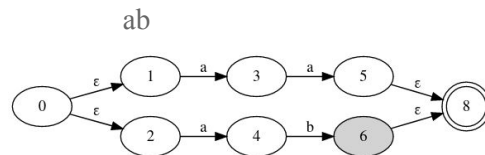
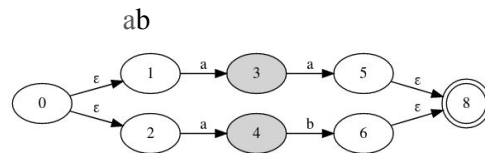
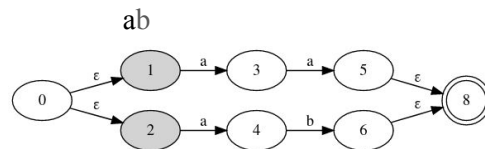
# Разбор с возвратами и без



fails, backtrack



matched successfully



matched successfully

# 1-однозначность

Под 1-однозначностью понимается свойство регулярных выражений, когда существует не более, чем один вариант успешного сопоставления для любого входного слова.

---

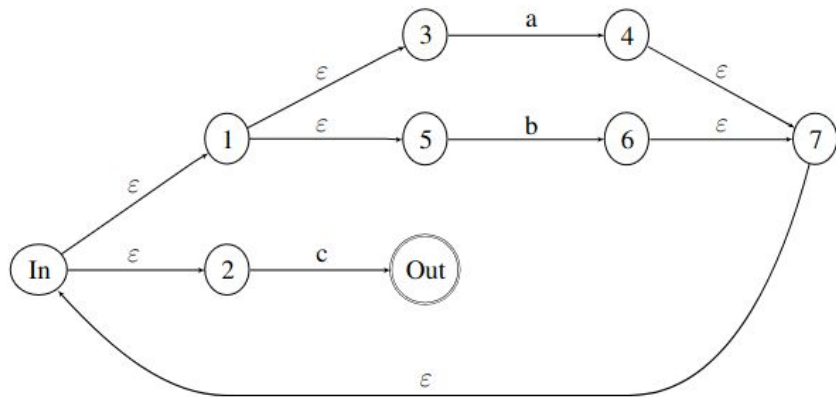
*Входную строку  $aaa$  по не 1-однозначному регулярному выражению  $(aa)^*a^*$  можно разобрать двумя разными способами:*

- *$aa$  соответствует  $(aa)^*$ , а соответствует  $a^*$ ;*
- *$aaa$  соответствует  $a^*$ , а шаблон  $(aa)^*$  сопоставляется с пустым словом.*

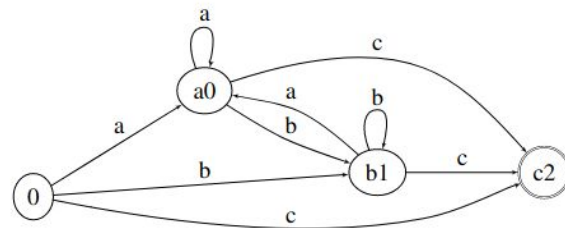
# Автомат Глушкова

$(a \mid b)^* c$

Автомат Томпсона



Автомат Глушкова



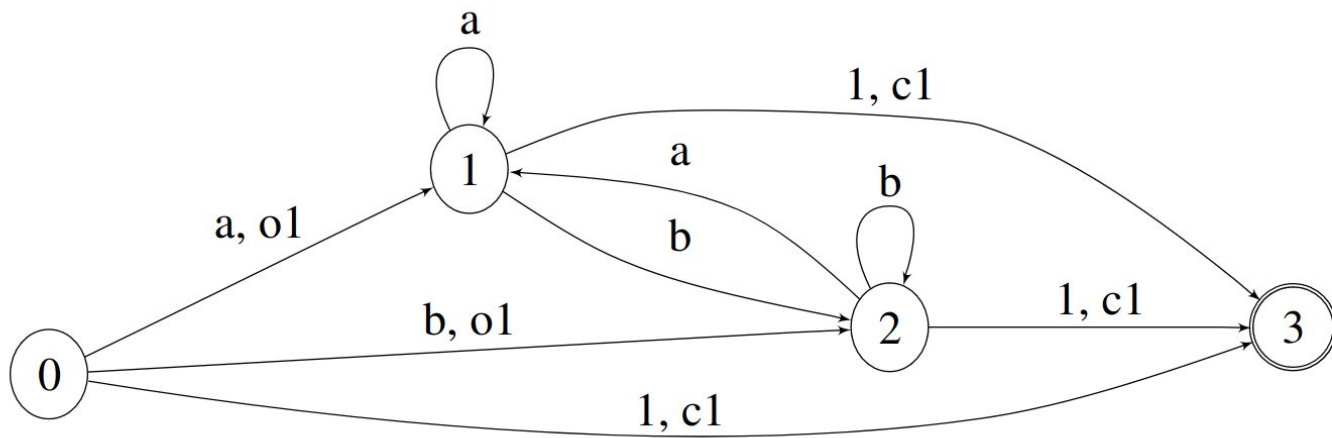
# MFA

Формально автомат с памятью, или MFA, определяется как пятерка элементов  $\langle Q, \Sigma, \delta, q_0, F \rangle$ , где  $Q$  — конечное множество состояний,  $\Sigma$  — это алфавит,  $q_0 \in Q$  — начальное состояние,  $F \subseteq Q$  — множество конечных состояний и  $\delta : Q \times (\Sigma \cup \{\varepsilon\} \cup \{1, 2, \dots, k\}) \rightarrow \mathcal{P}(Q \times \{o, c, \diamond\}^k)$  это функция переходов. Элементы  $o, c, \diamond$  называются инструкциями над памятью ( $o$  — открытие памяти,  $c$  — закрытие памяти,  $\diamond$  — сохранение памяти в состоянии, в котором она была).

*\*автомат с памятью считается  
детерминированным, если он  
детерминирован в классическом  
смысле над алфавитом и списком  
переменных*

# MFA

$$\underline{[_1a^* \mid b^*]_1 \& 1}$$





# Обращение классических регулярных выражений

- $\beta = \varepsilon \Rightarrow \text{reverse}(\beta) = \varepsilon;$
- $\beta = a \Rightarrow \text{reverse}(\beta) = a;$
- $\beta = A|B \Rightarrow \text{reverse}(\beta) = (\text{reverse}(A)|\text{reverse}(B));$
- $\beta = AB \Rightarrow \text{reverse}(\beta) = \text{reverse}(B)\text{reverse}(A);$
- $\beta = A^* \Rightarrow \text{reverse}(\beta) = (\text{reverse}(A))^*.$

---

*Регулярное выражение  $(a|b)^*a(a|b)$  недетерминировано. Однако при обращении получается детерминированное выражение  $(a|b)a(a|b)^*$ .*

# Ациклические регулярные выражения с обратными ссылками

- $\varepsilon$ ,  $\emptyset$ , а также все буквы алфавита  $\Sigma$  принадлежат  $\mathcal{P}_{\text{ACREG}}$ ;
- если  $\tau_1, \tau_2 \in \mathcal{P}_{\text{ACREG}}$ , тогда  $\tau_1\tau_2$ ,  $\tau_1 \mid \tau_2$ , а также  $\tau_1^*$  принадлежат  $\mathcal{P}_{\text{ACREG}}$ ;
- если  $i \in \mathbb{N}$ ,  $\tau \in \mathcal{P}_{\text{ACREG}}$ , и  $\tau$  не содержит  $[i\tau']_i$  ни для какого  $\tau'$ , тогда  $\&i$  и  $[i\tau]_i$  принадлежат  $\mathcal{P}_{\text{ACREG}}$ .
- $j \propto_r i$  — если в  $r$  встречается хотя бы одно подвыражение  $[i\tau]_i$  такое, что  $\&j$  входит в  $\tau$ ; транзитивное замыкание отношения  $\propto_r$  антирефлексивно.
- Для всякого  $\&i$ , входящего в  $r$  найдётся такое  $r'$ , что  $[ir']_i$  входит в  $r$ .

# ACREG != Алгебра Клини

$$\forall x, y, z \in \mathcal{A} (xy = yz \Rightarrow x^*y = yz^*)$$

- В  $\varepsilon$ -семантике положим  $x = [{}_1a]_1$ ,  $y = \&1$ ,  $z = aa$ . Тогда  $xy$  и  $yz$  оба задают язык  $\{aa\}$ , но  $\mathcal{L}(x^*y) = \{\varepsilon, a^{n+2}\}$ ,  $\mathcal{L}(yz^*) = \{a^{2 \cdot n}\}$  ( $n \geq 0$ ).
- В  $\emptyset$ -семантике положим  $x = [{}_1ba]_1$ ,  $y = \&1 \mid b$ ,  $z = ab \mid aba$ . Тогда  $xy$  и  $yz$  оба задают язык  $\{bab, baba\}$ , но  $\mathcal{L}(x^*y) = \{b, ba^{n+1}b, ba^{n+2}\}$ ,  $\mathcal{L}(yz^*) = \{b(ab \mid aba)^n\}$  ( $n \geq 0$ ).

# Неоднозначные чтения

---

*В регулярном выражении  $([a^*]_1|[b^*]_1|c)&1$  ссылка  $&1$  может относиться как к либо первому либо второму вариантам альтернативы, и соответственно соответствовать  $a^*$  либо  $b^*$ , так и к третьему варианту и соответствовать пустому слову.*

---

# lastinit-неоднозначность

Множество  $\text{last}_{i:\text{init}}(r)$  — это множество возможных выражений, инициализирующих ячейку памяти  $i$  после чтения выражения  $r$ .

Степень неоднозначности регулярного выражения  $r$  с обратными ссылками по переменной  $i$  —  $N_{\text{ambi}}r$ , равна  $n$ , если для  $k$ -го оператора чтения  $\&i$  в  $r$  и предшествующего ему выражения  $r'$   $|\text{last}_{i:\text{init}}r'| = n_k$ , и  $\sum_k n_k = n$ .

Степень  $\text{last}_{i:\text{init}}$ -неоднозначности регулярного выражения  $r$  с обратными ссылками —  $N_{\text{amb}}r$ , равна  $n$ , если  $n = \prod_i N_{\text{ambi}}r$ , где  $\&i$  входит в  $r$ .

---

Пусть  $r = [{}_1ba^*]_1[{}_2ca^*]_2([{}_1\&2ab]_1 \mid [{}_2bb^*]_2)^*$ . Тогда  $\text{last}_{2:\text{init}}r = \{ca^*, bb^*\}$ .

# СНФ

Ациклическое регулярное выражение  $r$  — в слабой ссылочной нормальной форме, если для каждого оператора чтения по ссылке  $\&i$  и предшествующего ему выражения  $r'$  множество  $\text{last}_{i:\text{init}} r'$  содержит единственный элемент.

Скажем, что  $r$  — в ссылочной нормальной форме (СНФ), если дополнительно к этому каждый оператор записи в память  $[_i r_0]_i$  инициализирует выражение  $r_0$ , которое входит в  $\text{last}_{i:\text{init}} r'$ , где  $r'$  предшествует оператору чтения  $\&i$ .

---

$$\begin{aligned} ([_1 a^*]_1 \mid [_1 b^*]_1 \mid c) \&1 &\rightarrow ([_1 a^*]_1 \&1 \mid [_1 b^*]_1 \&1 \mid c) \\ (c(\&1 \mid [_1 a^*]_1 \&1)b)^* &\rightarrow (cb)^* (c[_1 a^*]_1 \&1 b(c \&1 b)^*)^* \end{aligned}$$



# Преобразование в СНФ

- $(r_1 \mid r_2)r_3 \rightarrow (r_1r_3 \mid r_2r_3), r_1(r_2 \mid r_3) \rightarrow (r_1r_2 \mid r_1r_3)$
- $r_0r_1^*r_2 \rightarrow (r_0r_2 \mid r_0r_1^*r_1r_2)$
- $(r_1 \mid r_2)^*r_3 \rightarrow (r_1^*r_2)^*r_1^*r_3$
- $[_i(a \mid b)]_i \rightarrow ([_ia]_i \mid [_ib]_i)$
- $(r_1r_2)^*r_1 \rightarrow r_1(r_2r_1)^*$

# Обращение ACREG в СНФ

- $\beta = [{}_iA]_i$  и переменная уже была инициализована при обращении  $\Rightarrow$   
 $\text{reverse}(\beta) = \&i$ ,
- $\beta = [{}_iA]_i$  и переменная не была инициализована при обращении  $\Rightarrow$   
 $\text{reverse}(\beta) = [{}_i\text{reverse}(A)]_i$  и переменная добавляется в список инициализированных,
- $\beta = \&i$  и переменная уже была инициализована при обращении  $\Rightarrow$   
 $\text{reverse}(\beta) = \&i$ ,
- $\beta = \&i$  и переменная не была инициализована при обращении  $\Rightarrow$   
 $\text{reverse}(\beta) = [{}_i\text{reverse}(A)]_i$  и переменная добавляется в список инициализированных.



# RW блоки

- $[_1b^*]_1(\&1[_1a^*]_1)^* \rightarrow (b^*|[_1b^*]_1\&1([_1a^*]_1\&1) * a^*)$
- $([_2b^*]_2\&1[_1a^*]_1\&2)^*$
- $[_1a^*]_1(\&1[_1a^*]_1)^* \rightarrow (a^*|[_1a^*]_1(\&1[_1a^*]_1) * \&1a^*)$
- $[_2a^*d]_2([_1a^*]_1b\&1[_2\&1c * c]_2)^*$

# Примеры

$\{a^*\}:1/b\&1$

BNF:  $(\{a^*\}:1\&1|b\&1)$

Reverse:  $(\{a^*\}:1\&1|\&1b)$

$\{a^*\}:1/\{b^*\}:1/c\&1$

BNF:  $(\{a^*\}:1\&1|\{b^*\}:1\&1|c\&1)$

Reverse:  $(\{a^*\}:1\&1|\{b^*\}:1\&1|\&1c)$

$\{a^*\}:1(\&1\{a^*\}:1)^*$

BNF:  $(a^*|\{a^*\}:1(\&1\{a^*\}:1)^*\&1a^*)$

Reverse:  $(a^*|a^*\{a^*\}:1(\&1\{a^*\}:1)^*\&1)$

$\{b^*\}:2\&1\{a^*\}:1\&2)^*$

BNF:  $(\epsilon|\{b^*\}:2(a^*\&2|(\{a^*\}:1\&2\{b^*\}:2\&1)^*\{a^*\}:1\&2\{b^*\}:2\&1a^*\&2))^*$

Reverse:  $(\epsilon|(\{b^*\}:2a^*|\{b^*\}:2a^*\{a^*\}:1\&2\{b^*\}:2\&1(\{a^*\}:1\&2\&2\&1)^*)\&2))^*$

$c(\&1/\{a^*\}:1\&1)b)^*$

BNF:  $(cb)^*(c\{a^*\}:1\&1b(c\&1b)^*)^*$

Reverse:  $((b\{a^*\}:1c)^*b\&1\&1c)^*(bc)^*$

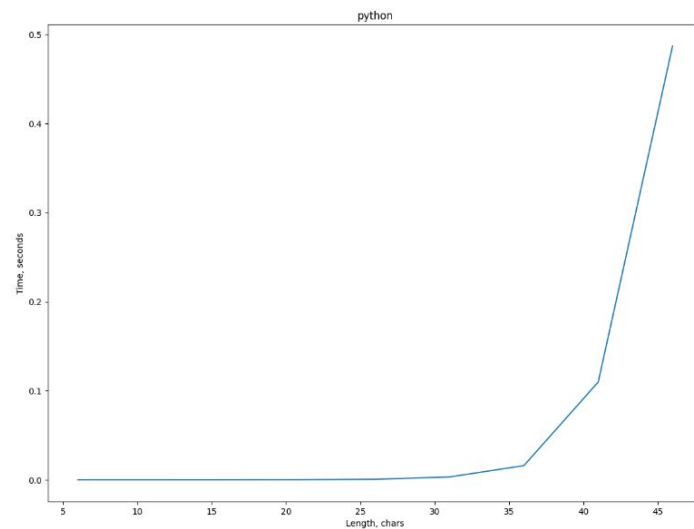
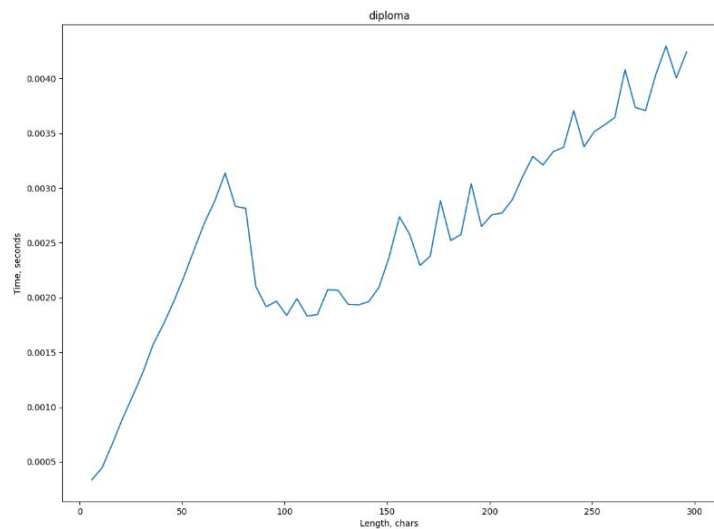
$((\{a^*\}:1/b)(\&1/b))^*$

BNF:  $(b|bb)^*(\{a^*\}:1\&1(bb)^*|a^*b(bb)^*|\{a^*\}:1\&1(bb)^*b\&1(bb)^*(b\&1(bb)^*)^*|\{a^*\}:1b(bb)^*b\&1(bb)^*(b\&1(bb)^*)^*)^*$

Reverse:  $((bb)^*\{a^*\}:1\&1|(bb)^*ba^*|((bb)^*\{a^*\}:1b)^*(bb)^*\&1b(bb)^*\&1\&1|((bb)^*\{a^*\}:1b)^*(bb)^*\&1b(bb)^*b\&1)^*(b|bb)^*$

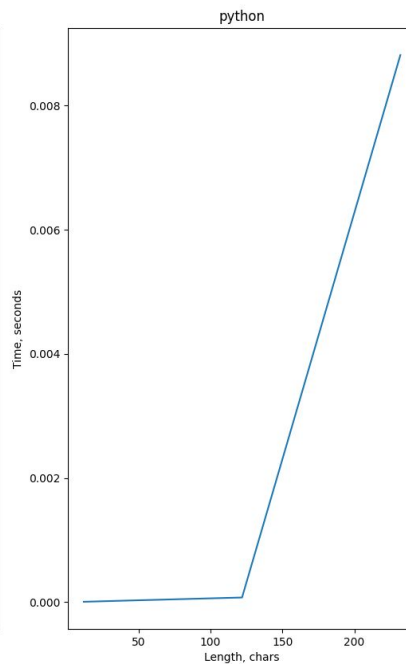
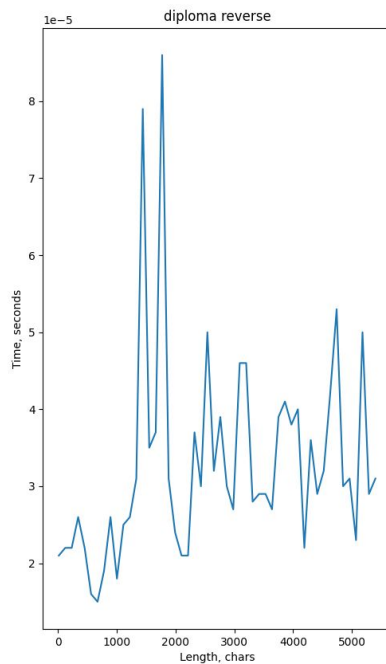
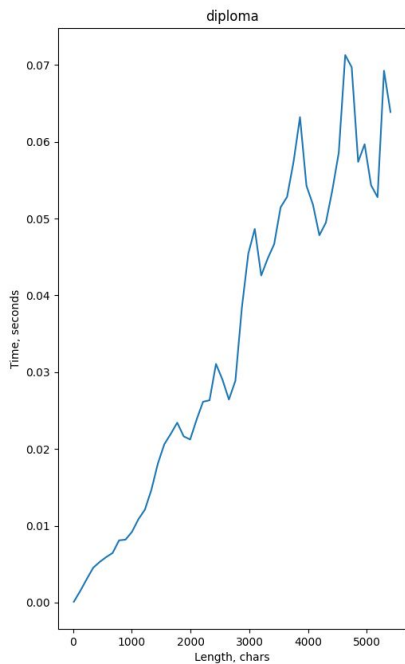
# Результаты

$\{a^*:1&1\}^*, ((a^*)^2)^*$



# Результаты

$(\{a^*\}:1b|&1)^*, ((a^*)b|2)^*$



# Список литературы

1. Kozen, Dexter. 1994. "A Completeness Theorem for Kleene Algebras and the Algebra of Regular Events." *Information and Computation* 110 (2): 366–90. <https://doi.org/https://doi.org/10.1006/inco.1994.1037>.
2. Campeanu, Cezar & Salomaa, Kai & Yu, Sheng. (2003). A Formal Study Of Practical Regular Expressions.. *Int. J. Found. Comput. Sci.* 14. 1007-1018. 10.1142/S012905410300214X.
3. Schmid, Markus. 2016. "Characterising REGEX Languages by Regular Languages Equipped with Factor-Referencing." *Information and Computation* (February). <https://doi.org/10.1016/j.ic.2016.02.003>.
4. Schmid, Markus. (2012). Inside the class of REGEX languages. *International Journal of Foundations of Computer Science*. 24. 73-84. 10.1007/978-3-642-31653-1\_8.
5. Brüggemann-Klein, Anne, and Derick Wood. 1998. "One-Unambiguous Regular Languages." *Information and Computation* 140 (2): 229–53. <https://doi.org/https://doi.org/10.1006/inco.1997.2688>.
6. Gruber, Hermann, and Stefan Gulan. 2010. "Simplifying Regular Expressions." In *Language and Automata Theory and Applications*, edited by Adrian-Horia Dediu, Henning Fernau, and Carlos Martín-Vide, 285–96. Berlin, Heidelberg: Springer Berlin Heidelberg.
7. Glushkov, V M. 1961. "THE ABSTRACT THEORY OF AUTOMATA." *Russian Mathematical Surveys* 16 (5): 1. <https://doi.org/10.1070/RM1961v016n05ABEH004112>.
8. Freydenberger, Dominik D., and Markus L. Schmid. 2018. "Deterministic Regular Expressions with Back-References." *CoRR abs/1802.01508*. <http://arxiv.org/abs/1802.01508>.
9. Hazel, Philip. n.d. "Официальное Руководство По PCRE2 (электронный ресурс)." <https://www.pcre.org/current/doc/html/index.html>.
10. Google "Репозиторий библиотеки Re2 (электронный ресурс)." <https://github.com/google/re2>.