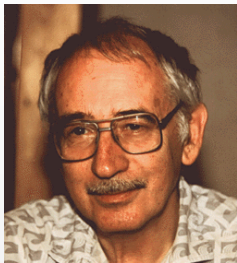


Краткое вступление



В.Ф. Турчин
(1931–2010)

- Функциональный язык Рефал
- Суперкомпиляция
- ...

-
- Семинары **МЕТА**: 2008–2016 в Переславле-Залесском
 - Приглашённые докладчики: Neil D. Jones, Simon Peyton-Jones

Распознаватели для формальных языков

Проблема слова

Пусть \mathcal{L} — формальный язык над алфавитом T , ω — терм в алфавите T . Проверить, верно ли, что $\omega \in \mathcal{L}$.

- Автоматы-распознаватели:
 - + Естественная реализация, эффективность
 - Неявные ограничения на свойства языка, слабая модифицируемость
- Алгебраические методы:
 - + Расширяемость (правила переписывания), семантика — база алгоритма
 - Неочевидная реализация, необходимость оптимизаций



Алгебраический язык как опора проектирования

Простейшие аналогии:

- ассоциативность операции — возможность асинхронного разбора
- коммутативность операции — возможность сортировки
- АСИ (ассоциативность, коммутативность, идемпотентность) — возможность представления в форме множества



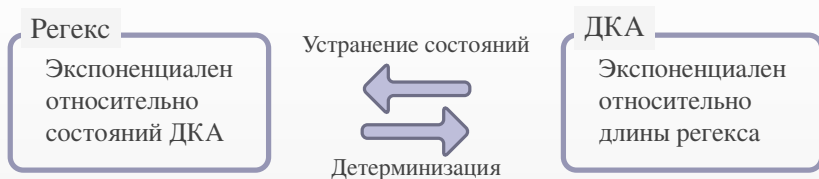
Регулярные языки

Академические регулярные выражения

- *Альтернатива (сложение, объединение) $|$ — ассоциативна, коммутативна, идемпотентна .*
- *Конкатенация (умножение) — ассоциативна.*
- *Унарная операция итерации $*$.*
- *Нет нетривиальных соотношений для нуля ($\underbrace{\emptyset}$).
пустой регекс*
- *Нетривиальные соотношения для единицы ($\underbrace{\varepsilon}$):
пустое слово*
 $\varepsilon \in x^*$, даже если $\varepsilon \notin x$.



Регулярные языки и их распознаватели



(Conway & Krob)

Устранение состояний сохраняет регекс по модулю правил:

- $x(yx)^* = (xy)^*x$
- $x^*(yx^*)^* = (x \mid y)^*$
- $\varepsilon \mid xx^* = x^*, \varepsilon \mid x^*x = x^*$
- $x(y \mid z) = xy \mid xz, (y \mid z)x = yx \mid zx$



Расширенные регексы — приоритеты и якоря

Альтернатива и конкатенация ассоциативны — **условно**, при переходе к именованным группам захвата.

Альтернатива коммутативна — **неверно**:

Регекс $((() | a)(a | b)^*$ не эквивалентен $(a | (())(a | b)^*$ из-за групп захвата.

- $\underbrace{\varepsilon}_{(())|a} \overbrace{aa}^{(a|b)^*}$ — группы 1 и 2 примут значение ε , группа примет значение a (последнее захваченное регексом $a|b$ значение).
- $\underbrace{a}_{(a|())} \overbrace{a}^{(a|b)^*}$ — группы 1 и 2 примут значение a , группа примет значение a .

Нет нетривиальных соотношений для нуля (\emptyset) — **неверно**:

Регекс $a \setminus A$ эквивалентен \emptyset : $\setminus A$ может лишь начинать выражение.



Производные Бржозовски

Множество $a^{-1}U = \{\omega \mid a\omega \in U\}$ называется производным Бржозовски множества U относительно a . Если $\varepsilon \in a^{-1}U$, тогда a распознаётся выражением U .

Λ_E положим равным $\{\varepsilon\}$, если $\varepsilon \in E$, и пустым множеством иначе.

- $a^{-1}\varepsilon = \emptyset, a^{-1}\emptyset = \emptyset$;
- $a^{-1}a = \{\varepsilon\}, a^{-1}b = \emptyset$;
- $a^{-1}(\Phi \mid \Psi) = a^{-1}(\Phi) \cup a^{-1}(\Psi)$;
- $a^{-1}(\Phi \Psi) = a^{-1}(\Phi)\Psi \cup \Lambda_\Phi a^{-1}(\Psi)$;
- $a^{-1}(\Phi^*) = a^{-1}(\Phi)\Phi^*$.

С помощью последовательного взятия производных можно свести задачу $\omega \in \mathcal{L}(R)$ к задаче $\varepsilon \in \omega^{-1}R$. Используя АСИ-эквивалентности, можно построить ДКА, помеченный производными языками, выполняющий такое распознавание.



Расширенные производные Бржозовски

Движок .NET:

- Отказ от перехода к автомату, производные используются непосредственно для разбора.
- Дополнительные условия для отслеживания единицы (синтаксически корректные якоря слов).
- Дополнительные условия для отслеживания нуля (синтаксически некорректные якоря слов).

Итог: алгоритм распознавания ω , линейный по длине ω .

- Возможные расширения на опережающие проверки.

D.Moseley et al: *Derivative Based Nonbacktracking Real-World Regex Matching with Backtracking Semantics*, Proceedings of the ACM on Programming Languages, Volume 7, Issue PLDI, Article No.: 148, pp. 1026–1049, 2023.



Расширенные регексы — обратные ссылки

Только регексы с именованными группами захвата.

Backref-регексы (слова со ссылками, Shmid):

$$\begin{cases} [k\tau]_k & (\text{захват в память}) \\ \&_k & (\text{чтение памяти}) \end{cases}$$

Пример: $[_1a^*]_1a^+b\&1$ определяет язык $\{a^m b a^n \mid m > n\}$

- ε -семантика (Shmid) — неинициализированная ссылка распознаёт $\{\varepsilon\}$;
- \emptyset -семантика (практические реализации) — неинициализированная ссылка распознаёт \emptyset .



Выбор не влияет на язык.



Проблема: циклическая память

- Выразительная сила и сложность анализа возрастает по сравнению с практическими регулярными выражениями.

Среди 3000 расширенных регексов со StackOverflow не нашлось ни одного циклического.

- Неочевидное синтаксическое свойство:
 - слово $([1\&2]_1[2a\&1]_2[1a^*]_1)^*$ циклическое,
 - слово $([1\&2]_1[1a^*]_1[2a\&1]_2)^*$ нециклическое.



Семантически корректная переименовка

	$[1a^*]_1$ $\mapsto [3a^*]_3$	$[1\&2]_1$ $\mapsto [3\&2]_3$	Циклы после подст.
$([1\&2]_1[2a\&1]_2[1a^*]_1)^*$	✓	✗	✓
$([1\&2]_1[1a^*]_1[2a\&1]_2)^*$	✗	✓	✗



Ациклические регексы

- Ограниченная переинициализация зависимых ячеек памяти
- Скобки групп захвата — операции над памятью
- Могут быть корректно переименованы в регексы без переинициализаций

ACREG — идемпотентное полукольцо, удовлетворяющее теоремам Конвея – Кроба



Конечный автомат с памятью

Пятерка $\langle Q, \Sigma, \delta, q_0, F \rangle$:

- Q — множество состояний автомата;
- Σ — входной (терминальный) алфавит автомата;
- δ — множество правил перехода. Переходы помечаются действиями над памятью: o — открытие ячейки, c — закрытие ячейки, r — сброс до ε ;
- $q_0 \in Q$ — начальное состояние, $F \subseteq Q$ — множество конечных состояний.

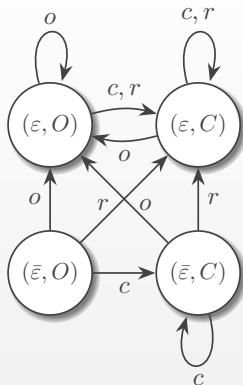
Операция сброса не применяется в оригинальной статье Шмида и необходима, чтобы сделать алгебру действий над памятью композиционно замкнутой.



MFA: состояния памяти

- текущее состояние;
- содержимое памяти;
- конфигурации ячеек:

$$\begin{cases} O & (\text{open}) \\ C & (\text{closed}) \end{cases}$$



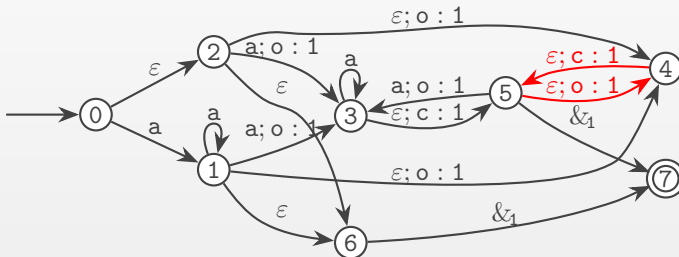
- Начальная конфигурация памяти: $(q_0, w, (\varepsilon, C), \dots, (\varepsilon, C))$.
- Действия над ячейками памяти происходят до чтения с ленты.



Необходимость откатов и пустые действия

Расширенное регулярное выражение: $a^*[a^*] : 1^* \&_1$

Без сбросов памяти необходимо возникают циклы по пустому слову.



Необходимость откатов и пустые действия

Расширенное регулярное выражение: $a^*[a^*] : 1^*\&_1$

Без переходов по пустому слову с операцией сброса распознавание эффективно.

