

Deriving (finite) intersection non-emptiness, courtesy of Brzozowski

Breandan Considine

March 11, 2025

1 Syntax and semantics

Definition 1 (Generalized Regex). Let E be an expression defined by the grammar:

$$E ::= \emptyset \mid \varepsilon \mid \Sigma \mid E \cdot E \mid E \vee E \mid E \wedge E$$

Semantically, we interpret these expressions as denoting regular languages:

$$\begin{array}{ll} \mathcal{L}(\emptyset) = \emptyset & \mathcal{L}(S \cdot T) = \mathcal{L}(S) \times \mathcal{L}(T) \\ \mathcal{L}(\varepsilon) = \{\varepsilon\} & \mathcal{L}(S \vee T) = \mathcal{L}(S) \cup \mathcal{L}(T) \\ \mathcal{L}(a) = \{a\} & \mathcal{L}(S \wedge T) = \mathcal{L}(S) \cap \mathcal{L}(T) \end{array}$$

Definition 2 (Brzozowski, 1964). To compute the quotient $\partial_a(L) = \{b \mid ab \in L\}$, we:

$$\begin{array}{ll} \partial_a(\emptyset) = \emptyset & \delta(\emptyset) = \emptyset \\ \partial_a(\varepsilon) = \emptyset & \delta(\varepsilon) = \varepsilon \\ \partial_a(b) = \begin{cases} \varepsilon & \text{if } a = b \\ \emptyset & \text{if } a \neq b \end{cases} & \delta(a) = \emptyset \\ \partial_a(S \cdot T) = (\partial_a S) \cdot T \vee \delta(S) \cdot \partial_a T & \delta(S \cdot T) = \delta(S) \wedge \delta(T) \\ \partial_a(S \vee T) = \partial_a S \vee \partial_a T & \delta(S \vee T) = \delta(S) \vee \delta(T) \\ \partial_a(S \wedge T) = \partial_a S \wedge \partial_a T & \delta(S \wedge T) = \delta(S) \wedge \delta(T) \end{array}$$

Theorem 1 (Recognition). For any regex R and $\sigma : \Sigma^*$, $\sigma \in \mathcal{L}(R) \iff \varepsilon \in \mathcal{L}(\partial_\sigma R)$, where:

$$\partial_\sigma(R) : RE \rightarrow RE = \begin{cases} R & \text{if } \sigma = \varepsilon \\ \partial_b(\partial_a R) & \text{if } R = a \cdot b, a \in \Sigma \end{cases}$$

Theorem 2 (Generation). For any (ε, \wedge) -free regex, R , to generate a witness $\sigma \sim \mathcal{L}(R)$:

$$\begin{array}{l} follow(R) : RE \rightarrow 2^\Sigma = \begin{cases} \{R\} & \text{if } R \in \Sigma \\ follow(S) & \text{if } R = S \cdot T \\ follow(S) \cup follow(T) & \text{if } R = S \vee T \end{cases} \\ choose(R) : RE \rightarrow \Sigma^+ = \begin{cases} R & \text{if } R \in \Sigma \\ (s \sim follow(R)) \cdot choose(\partial_s R) & \text{if } R = S \cdot T \\ choose(R' \sim \{S, T\}) & \text{if } R = S \vee T \end{cases} \end{array}$$

2 Language intersection

Theorem 3 (Bar-Hillel, 1961). *For any context-free grammar (CFG), $G = \langle V, \Sigma, P, S \rangle$, and nondeterministic finite automata, $A = \langle Q, \Sigma, \delta, I, F \rangle$, there exists a CFG $G_\cap = \langle V_\cap, \Sigma_\cap, P_\cap, S_\cap \rangle$ such that $\mathcal{L}(G_\cap) = \mathcal{L}(G) \cap \mathcal{L}(A)$.*

Definition 3 (Salomaa, 1973). One could construct G_\cap like so,

$$\frac{q \in I \quad r \in F}{(S \rightarrow qSr) \in P_\cap} \vee \frac{(A \rightarrow a) \in P \quad (q \xrightarrow{a} r) \in \delta}{(qAr \rightarrow a) \in P_\cap} \uparrow \frac{(w \rightarrow xz) \in P \quad p, q, r \in Q}{(pwr \rightarrow (pxq)(qzr)) \in P_\cap} \bowtie$$

however most synthetic productions in P_\cap will be non-generating or unreachable.

Theorem 4 (Considine, 2025). *For every CFG, G , and every acyclic NFA (ANFA), A , there exists a decision procedure $\varphi : \text{CFG} \rightarrow \text{ANFA} \rightarrow \mathbb{B}$ such that $\varphi(G, A) \models [\mathcal{L}(G) \cap \mathcal{L}(A) \neq \emptyset]$ which requires $\mathcal{O}((\log |Q|)^c)$ time using $\mathcal{O}((|V||Q|)^k)$ parallel processors for some $c, k < \infty$.*

Proof sketch. WTS there exists a path $p \rightsquigarrow r$ in A such that $p \in I, r \in F$ where $S \vdash p \rightsquigarrow r$.

There are two cases, at least one of which must hold for $w \in V$ to parse a given $p \rightsquigarrow r$ pair:

1. p steps directly to r in which case it suffices to check $\exists s. ((p \xrightarrow{s} r) \in \delta \wedge (w \rightarrow s) \in P)$, or,
2. there is some midpoint $q \in Q$, $p \rightsquigarrow q \rightsquigarrow r$ such that $\exists x, z. ((w \rightarrow xz) \in P \wedge \underbrace{p \rightsquigarrow q}_x \underbrace{q \rightsquigarrow r}_z)$.

This decomposition suggests a dynamic programming solution. Let M be a matrix of type $RE^{|Q| \times |Q| \times |V|}$ indexed by Q . Since we assumed δ is acyclic, there exists a topological sort of δ imposing a total order on Q such that M is strictly upper triangular (SUT). Initiate it thusly:

$$M_0[r, c, v] = \bigvee_{a \in \Sigma} \{a \mid (v \rightarrow a) \in P \wedge (q_r \xrightarrow{a} q_c) \in \delta\} \quad (1)$$

The algebraic operations $\oplus, \otimes : RE^{2|V|} \rightarrow RE^{|V|}$ will be defined elementwise:

$$[\ell \oplus r]_w = [\ell_w \vee r_w] \quad (2)$$

$$[\ell \otimes r]_w = \bigvee_{x, z \in V} \{\ell_x \cdot r_z \mid (w \rightarrow xz) \in P\} \quad (3)$$

By slight abuse of notation, we will redefine the matrix exponential over this domain as:

$$\exp(M) = \sum_{i=0}^{\infty} M_0^i = \sum_{i=0}^{|Q|} M_0^i \text{ (since } M \text{ is SUT.)} \quad (4)$$

To solve for the fixpoint, we can instead use exponentiation by squaring:

$$S(2n) = \begin{cases} M_0, & \text{if } n = 1, \\ S(n) + S(n)^2 & \text{otherwise.} \end{cases} \quad (5)$$

Therefor, we only need a maximum of $\lceil \log_2 |Q| \rceil$ sequential steps to reach the fixpoint. Finally,

$$S_\cap = \bigvee_{q \in I, q' \in F} \exp(M)[q, q', S] \text{ and } \varphi = [S_\cap \neq \emptyset] \quad (6)$$

To decode a witness in case of non-emptiness, we simply **choose**(φ). \square

3 Combinatorics

To enumerate, we first need $|\mathcal{L}(R)|$, which is denoted $|R|$ for brevity.

Definition 4 (Cardinality). $|R| : RE \rightarrow \mathbb{N} = \begin{cases} 1 & \text{if } R \in \Sigma \\ S \times T & \text{if } R = S \cdot T \\ S + T & \text{if } R = S \vee T \end{cases}$

Theorem 5 (Enumeration). *To enumerate, invoke $\bigcup_{i=0}^{|R|} \text{enum}(R, i)$:*

$$\text{enum}(R, n) : RE \times \mathbb{N} \rightarrow \Sigma^* = \begin{cases} R & \text{if } R \in \Sigma \\ \text{enum}(S, \lfloor \frac{n}{|T|} \rfloor) \cdot \text{enum}(T, n \bmod |T|) & \text{if } R = S \cdot T \\ (n < |S|) ? \text{enum}(S, n) : \text{enum}(T, n - |S|) & \text{if } R = S \vee T \end{cases}$$

4 Fermi estimation

Suppose a PRAM-based model with 1 TFLOP throughput $\approx 2^{40}$ FLOPS and assume a 2^8 FLOP overhead per entry, n^3 entries per matrix multiplication, and approximately $\log_2(n)$ matrix multiplications. Thus, we would expect somewhere on the order of $2^{10} n^3 \log_2(n)$ total FLOPS per language intersection instance with a ballpark latency of:

$$T \approx \frac{2^8 n^3 \log_2(n)}{2^{40}} = \frac{n^3 \log_2(n)}{2^{32}} \text{ s}$$

For example assuming, $n^3 = (2^9)^3 = 2^{27}$, $\log_2(n) = 9$, hence we can estimate:

$$T \approx \frac{2^{27} \cdot 9}{2^{32}} = \frac{9}{2^5} = \frac{9}{32} \approx 280 \text{ ms}$$

Further optimizations are likely realizable via bitpacking and vectorization.

5 Future work

Broadly interested in questions related to formal languages and finite model theory, following the Carnapian program of logical syntax. Encoding semantics into syntax would allow us to do type checking and static analysis in the parser. A few lines of attack here:

Result 1 (Büchi–Elgot–Trakhtenbrot). A language is regular iff it is MSO-definable. For every MSO formula, there is a corresponding FSA. Complexity may be superexponential.

Result 2 (Pentus, 1993). Lambek categorial grammars, a weak kind of substructural logic, recognize exactly the context-free languages.

Result 3 (DeYoung & Pfenning, 2016). Describes a certain equivalence between subsingleton logic, a weak kind of linear logic, and automata.

Result 4 (Knuth & Wegner, 1968). Attribute grammars permit some notation of semanticity.

6 Acknowledgements

The author wishes to thank Fr. Paul Pomkowski for his phonetic advice.