

Syntax Repair as Language Intersection

ANONYMOUS AUTHOR(S)

We introduce a new technique for correcting syntax errors in arbitrary context-free languages. Our work addresses the problem of syntax error correction, which we solve by defining a finite language that provably generates every repair within a certain edit distance. To do this, we adapt the Bar-Hillel construction from formal languages, guaranteeing the language is sound and complete with respect to the programming language grammar. This technique also admits a polylogarithmic time algorithm for deciding intersection nonemptiness between CFLs and acyclic NFAs, the first of its kind in the parsing literature.

1 INTRODUCTION

During programming, one invariably encounters a recurring scenario in which the editor occupies an unparseable state, either due to an unfinished or malformed piece of code. Faced with this predicament, programmers must spend time to locate and repair the error before proceeding. We attempt to solve this problem automatically by generating a list of candidate repairs which contains with high probability the true repair, assuming this repair differs by no more than a few edits.

2 BACKGROUND

We will first give some background on Brzozowski differentiation. Let us consider the generalized regular expression (GRE) fragment containing concatenation, conjunction and disjunction:

Definition 2.1 (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 can interpret these expressions as denoting regular languages:

$$\begin{array}{ll} \mathcal{L}(\emptyset) = \emptyset & \mathcal{L}(x \cdot z) = \mathcal{L}(x) \times \mathcal{L}(z)^1 \\ \mathcal{L}(\varepsilon) = \{\varepsilon\} & \mathcal{L}(x \vee z) = \mathcal{L}(x) \cup \mathcal{L}(z) \\ \mathcal{L}(a) = \{a\} & \mathcal{L}(x \wedge z) = \mathcal{L}(x) \cap \mathcal{L}(z) \end{array}$$

Brzozowski introduces the concept of differentiation, which allows us to quotient a regular language by some given prefix.

Definition 2.2 (Brzozowski, 1964). To compute the quotient $\partial_a(L) = \{b \mid ab \in L\}$, we:

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

Primarily, this gadget was designed to handle membership queries, for which purpose it has received considerable attention in recent years:

¹Or $\{a \cdot b \mid a \in \mathcal{L}(x) \wedge b \in \mathcal{L}(z)\}$ to be more precise, however we make no distinction.

THEOREM 2.3 (RECOGNITION). For any regex e and $\sigma : \Sigma^*$, $\sigma \in \mathcal{L}(e) \iff \varepsilon \in \mathcal{L}(\partial_\sigma e)$, where:

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

It can also be used, however, to decode witnesses. We will first focus on the nonempty disjunctive fragment, and define this process in two steps:

THEOREM 2.4 (GENERATION). For any nonempty (ε, \wedge) -free regex, e , to witness $\sigma \in \mathcal{L}(e)$:

$$\begin{aligned} \text{follow}(e) : E \rightarrow 2^\Sigma &= \begin{cases} \{e\} & \text{if } e \in \Sigma \\ \text{follow}(x) & \text{if } e = x \cdot z \\ \text{follow}(x) \cup \text{follow}(z) & \text{if } e = x \vee z \end{cases} \\ \text{choose}(e) : E \rightarrow \Sigma^+ &= \begin{cases} e & \text{if } e \in \Sigma \\ (s \overset{\$}{\leftarrow} \text{follow}(e)) \cdot \text{choose}(\partial_s e) & \text{if } e = x \cdot z \\ \text{choose}(e' \overset{\$}{\leftarrow} \{x, z\}) & \text{if } e = x \vee z \end{cases} \end{aligned}$$

Here, we use the $\overset{\$}{\leftarrow}$ operator to denote probabilistic choice, however any deterministic choice function will also suffice to generate a witness. Now we are equipped to handle conjunction.

2.1 Language intersection

THEOREM 2.5 (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 2.6 (Salomaa, 1973). One could construct G_\cap like so,

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

However most synthetic productions in P_\cap will be non-generating or unreachable. This naïve method will construct a synthetic production for state pairs which are not even connected by any path, which is clearly excessive. We will instead proceed by constructing a parse chart that represents the intersection.

THEOREM 2.7. 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 $p \rightsquigarrow r \vdash S$.

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 $\exists a. ((p \overset{a}{\rightarrow} r) \in \delta \wedge (w \rightarrow a) \in P)$, or,
- (2) there is some midpoint $q \in Q$, $p \rightsquigarrow q \rightsquigarrow r$ such that $\exists x, z. ((w \rightarrow xz) \in P \wedge \overbrace{p \rightsquigarrow q}^w, \overbrace{q \rightsquigarrow r}^z)$.

This decomposition immediately suggests a dynamic programming solution. Let M be a matrix of type $E^{|Q| \times |Q| \times |V|}$ 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:

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

The algebraic operations $\oplus, \otimes : E^{2|V|} \rightarrow E^{|V|}$ 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 notation², 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 (S_{\cap}) . □

2.2 Combinatorics

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

$$\text{Definition 2.8 (Cardinality). } |e| : E \rightarrow \mathbb{N} = \begin{cases} 1 & \text{if } R \in \Sigma \\ x \times z & \text{if } e = x \cdot z \\ x + z & \text{if } e = x \vee z \end{cases}$$

THEOREM 2.9 (ENUMERATION). To enumerate, invoke $\bigcup_{i=0}^{|R|} \{enum(R, i)\}$:

$$enum(e, n) : E \times \mathbb{N} \rightarrow \Sigma^* = \begin{cases} e & \text{if } R \in \Sigma \\ enum(x, \lfloor \frac{n}{|z|} \rfloor) \cdot enum(z, n \bmod |z|) & \text{if } e = x \cdot z \\ enum((x, z)_{\min(1, \lfloor \frac{n}{|x|} \rfloor)}, n - |x| \min(1, \lfloor \frac{n}{|x|} \rfloor)) & \text{if } e = x \vee z \end{cases}$$

²Traditionally, there is a $\frac{1}{k!}$ factor to suppress exploding entries.

3 MEASURING THE LANGUAGE INTERSECTION

We will now attempt to put a probability distribution over the language intersection. We will start with a few cursory but illuminative approaches, then proceed towards a more refined solution.

3.1 Exact enumeration

A brute force solution would be to generate every path and rank every one by its probability. It should be obvious why is unviable due its worst case complexity, but bears mentioning due to its global optimality. In certain cases, it can be realized when the intersection language is small.

3.2 Mode collapse

Ordinarily, we would use top-down PCFG sampling, however in the case of non-recursive CFGs, this method is highly degenerate, exhibiting poor sample diversity. Consider an illustrative pathological case for top-down ancestral (TDA) sampling:

$$\begin{aligned} S &\rightarrow AB \ (0.9999) & S &\rightarrow CC \ (0.0001) \\ A &\rightarrow a \ (1) & B &\rightarrow b \ (1) & C &\rightarrow a \left(\frac{1}{26}\right) \mid \dots \mid z \left(\frac{1}{26}\right) \end{aligned}$$

TDA sampling will almost always generate the string ab , but most of the language is concealed in the hidden branch, $S \rightarrow CC$. Although contrived example, it illustrates precisely why TDA sampling is unviable: we want a sampler that matches the true distribution over the finite CFL, not the PCFG's local approximation thereof.

3.3 Ambiguity

Another approach would be to sample trees and rerank them by their PCFG score. More pernicious is the issue of ambiguity. Since the CFG can be ambiguous, this causes certain repairs to be overrepresented, resulting in a subtle bias. Consider for example,

LEMMA 3.1. *If the FSA, α , is ambiguous, then the intersection grammar, G_{\cap} , can be ambiguous.*

PROOF. Let ℓ be the language defined by $G = \{S \rightarrow LR, L \rightarrow (, R \rightarrow)\}$, where $\alpha = L(\underline{\sigma}, 2)$, the broken string σ is $) ($, and $\mathcal{L}(G_{\cap}) = \ell \cap \mathcal{L}(\alpha)$. Then, $\mathcal{L}(G_{\cap})$ contains the following two identical repairs: $) ($ with the parse $S \rightarrow q_{00}Lq_{21} q_{21}Rq_{22}$, and $) ($ with the parse $S \rightarrow q_{00}Lq_{11} q_{11}Rq_{22}$. \square

We would like the underlying sample space to be a proper set, *not* a multiset.