

# 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 this language is sound and complete with respect to a programming language's 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

When programming, one invariably encounters a recurring scenario in which the editor occupies an unparseable state. Faced with this predicament, programmers must spend time to locate and repair the error before proceeding. We 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.

Prior research in syntax repair uses stochastic and deterministic methods. In the case of neural program repair, the model is typically used as a base distribution, however most neural language models are only approximate inference methods and thus require some form of postprocessing and rejection sampling to ensure the results conform to the grammar. In contrast, most rule-based syntax repair methods, while exact, are too weak to fully model the full distribution of natural source code and must rely on local heuristics.

The primary issue with both of these methods is that they generate a very low number of repairs. Even if the repair model is highly accurate, it is unlikely to generate the true repair when there are a large number of approximately equidistant or equiprobable repairs.

To address this problem, our method explicitly generates and ranks every repair within a fixed edit distance. This ensures that the true repair is retrieved and scored, as long as it is only a few edits apart. Specifically, we employ the Bar-Hillel construction from formal language theory to generate a regular expression representing the finite language, then enumerate up to  $2^{16}$  distinct words from this language, rerank the enumerated set using a neural language model, and present the top dozen or so results to the user, all within a few hundred milliseconds.

We evaluate our approach on a dataset of human syntax errors and fixes fewer than five lexical edits apart and shorter than 120 tokens, large enough to fit a few lines of source code in realistic programming languages and editing scenarios. Our work shows this technique is extremely effective at predicting the true repair across a dataset of Python source code, on average 5x more accurately than the previous state of the art at the same latency.

## 2 BACKGROUND

Recall that a CFG,  $\mathcal{G} = \langle \Sigma, V, P, S \rangle$ , is a quadruple consisting of terminals ( $\Sigma$ ), nonterminals ( $V$ ), productions ( $P: V \rightarrow (V \mid \Sigma)^+$ ), and a start symbol, ( $S$ ). Every CFG is reducible to so-called *Chomsky Normal Form* [2],  $P': V \rightarrow (V^2 \mid \Sigma)$ , where every production is either (1) a binary production  $w \rightarrow xz$ , or (2) a unit production  $w \rightarrow t$ , where  $w, x, z : V$  and  $t : \Sigma$ . For example:

$$G = \{ S \rightarrow SS \mid (S) \mid ( ) \} \implies G' = \{ S \rightarrow QR \mid SS \mid LR, \quad R \rightarrow ), \quad L \rightarrow (, \quad Q \rightarrow LS \}$$

Likewise, a finite state automaton (FSA) is a quintuple  $\mathcal{A} = \langle Q, \Sigma, \delta, I, F \rangle$ , where  $Q$  is a finite set of states,  $\Sigma$  is a finite alphabet,  $\delta \subseteq Q \times \Sigma \times Q$  is the transition function, and  $I, F \subseteq Q$  are the set of initial and final states, respectively. We will adhere to this notation in the following sections.

There is an equivalent characterization of the regular languages using an inductively defined datatype which is often more elegant to work with. 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:

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

Brzowski [1] introduces the concept of differentiation, which allows us to quotient a regular language by some given prefix.

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

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

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

**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:

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

Variations on this basic procedure can also be used for functional parsing and regular expression tasks. Brzowski's derivative can also be used to decode witnesses. We will first focus on the nonempty disjunctive fragment, and define this process in two steps:

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

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

$$\text{follow}(e) : E \rightarrow 2^\Sigma = \begin{cases} \{e\} & \text{if } e \in \Sigma \\ \text{follow}(x) & \text{if } e = x \cdot z \\ \text{follow}(x) \cup \text{follow}(z) & \text{if } e = x \vee z \end{cases}$$

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

Here, we use the  $\xleftarrow{\$}$  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

We will now define intersection in a slightly more expressive manner, which has the added benefit of being more readily parallelizable. Recall every regular language is context-free. Therefore, to take the intersection between two regular languages, we can treat one as a CFL, which often admits a more compact representation. Alternatively, we can take the intersection between a truly non-regular CFL (such as a programming language syntax) and some regular language.

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  $\mathcal{L}(G_\cap) = \mathcal{L}(G) \cap \mathcal{L}(A)$ .*

Salomaa [5] introduces a direct, but inefficient construction:

Definition 2.6 (Salomaa, 1973). One could construct  $G_\cap$  like so,

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

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.

### 3 INFORMAL STATEMENT

Assume there exists a transducer from Unicode tokens to grammatical tokens,  $\tau : \Sigma_U^* \rightarrow \Sigma_G^*$ . In the compiler nomenclature  $\tau$  is called a *lexer* and would typically be regular under mild assumptions. Lexing abstracts over named identifiers and primitively typed literals, and removes whitespace or otherwise replaces them with an alphanumeric identifier in whitespace-sensitive languages. In this paper, we do not consider  $\tau$  and strictly deal with languages over  $\Sigma_G^*$ , or simply  $\Sigma^*$  for brevity.

Now suppose we have a syntax,  $\ell \subset \Sigma^*$ , containing every acceptable program. A syntax error is an unacceptable string,  $\sigma \notin \ell$ , that we wish to repair. We can model syntax repair as a language intersection between a context-free language (CFL) and a regular language. Henceforth,  $\sigma$  will always and only be used to denote a syntactically invalid string whose target language is known.

Given the source code for a computer program  $\sigma$  and a grammar  $G$ , our goal is to find every valid string  $\tilde{\sigma}$  consistent with the grammar  $G$  and within a certain edit distance,  $d$ . Consider the language of valid strings within a given Levenshtein distance from a reference string  $\sigma$ . We can intersect the language given by the Levenshtein automaton with the language of all valid programs given by the grammar  $G$ . The resulting language  $\mathcal{L}(G_\cap)$  will contain every repair within the designated edit distance.

Once the repairs are retrieved, they can be mapped back into raw source code tokens in a straightforward manner. The core algorithm is formatting-oblivious. After the repairs are returned, we compute an edit alignment over lexical tokens, restore unmodified tokens to their Unicode form, add placeholders for fresh names, numbers, and string literals, then finally apply an off-the-shelf code formatter to return the concrete suggestions. Both the preprocessing and the formatting steps are tangential to this paper, in which we confine ourselves to a lexical alphabet.

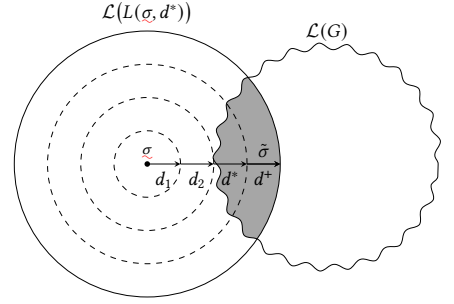


Fig. 1. CFL intersection with the local edit region of a given broken code snippet.

#### 4 FORMAL STATEMENT

*Definition 4.1 (Bounded Levenshtein-CFL reachability).* Given a CFL,  $\ell$ , and an invalid string,  $\underline{\sigma} : \bar{\ell}$ , find every valid string reachable within  $d$  edits of  $\underline{\sigma}$ , i.e., letting  $\Delta$  be the Levenshtein metric and  $L(\underline{\sigma}, d) = \{\sigma' \mid \Delta(\underline{\sigma}, \sigma') \leq d\}$  be the Levenshtein  $d$ -ball, we seek to find  $\ell_{\cap} = L(\underline{\sigma}, d) \cap \ell$ .

As the admissible set  $\ell_{\cap}$  is typically under-constrained, we want a procedure which surfaces natural and valid repairs over unnatural but valid repairs:

*Definition 4.2 (Ranked repair).* Given a finite language  $\ell_{\cap} = L(\underline{\sigma}, d) \cap \ell$  and a probabilistic language model  $P_{\theta} : \Sigma^* \rightarrow [0, 1] \subset \mathbb{R}$ , the ranked repair problem is to find the top- $k$  maximum likelihood repairs under the language model. That is,

$$R(A, P_{\theta}) = \operatorname{argmax}_{\sigma \in \ell_{\cap}, |\sigma| \leq k} \sum_{\sigma \in \sigma} P_{\theta}(\sigma) \quad (1)$$

A popular approach to ranked repair involves learning a distribution over strings, however this is highly sample-inefficient and generalizes poorly to new languages. Approximating a distribution over  $\Sigma^*$  forces the model to jointly learn syntax and stylometry. Furthermore, even with an extremely efficient approximate sampler for  $\sigma \sim \ell_{\cap}$ , due to the size of  $\ell$  and  $L(\underline{\sigma}, d)$ , it would be intractable to sample either  $\ell$  or  $L(\underline{\sigma}, d)$ , reject duplicates, then reject invalid ( $\sigma \notin \ell$ ) or unreachable ( $\sigma \notin L(\underline{\sigma}, d)$ ) edits, and completely out of the question to sample  $\sigma \sim \Sigma^*$  as do many neural language models.

As we will demonstrate, the ranked repair problem can be factorized into a bilevel objective: first maximal retrieval, then ranking. Instead of working with strings, we will explicitly construct a grammar which soundly and completely generates the set  $\ell \cap L(\underline{\sigma}, d)$ , then retrieve repairs from its language. By ensuring retrieval is sufficiently precise and exhaustive, maximizing likelihood over the retrieved set can be achieved with a much simpler, syntax-oblivious language model.

Assuming we have a grammar that recognizes the Levenshtein-CFL intersection, the question then becomes how to maximize the number of unique valid sentences in a given number of samples. Top-down incremental sampling with replacement eventually converges to the language, but does so superlinearly [?]. Due to practical considerations including latency, we require the sampler to converge linearly, ensuring with much higher probability that natural repairs are retrieved in a timely manner. This motivates the need for a specialized generating function. More precisely,

*Definition 4.3 (Linear convergence).* Given a finite CFL,  $\ell$ , we want a randomized generating function,  $\phi : \mathbb{N}_{\leq |\ell|} \rightarrow 2^{\ell}$ , whose rate of convergence is linear in expectation, i.e.,  $\mathbb{E}_{i \in [1, n]} |\phi(i)| \propto n$ .

This will ensure that if  $|\ell_{\cap}|$  is sufficiently small and enough samples are drawn,  $\phi$  is sure to include a representative subset, and additionally, will terminate after exhausting all valid repairs.

To satisfy Def. 4.3, we can construct a bijection from syntax trees to integers (§ ??), sample integers uniformly without replacement, then decode them as trees. This will produce a set of unique trees, and each tree, assuming grammatical unambiguity, will correspond to a unique sentence in the language. Finally, sentences can be scored and ranked by likelihood under a language model.

Otherwise, if the grammar,  $G_{\ell}$ , is ambiguous, it can be translated into a DFA, then decoded (§ ??) using an autoregressive language model or any suitably fast scoring function of the implementer's choice. In our case, we use a low-order Markov model for its inference speed, data efficiency, and simplicity. So long as the decoder samples  $\ell$  without replacement, it will satisfy Def. 4.3.

## 5 METHOD

Our method is to treat finite language intersections as matrix exponentiation.

**THEOREM 5.1.** *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.** 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  $\exists a. ((p \xrightarrow{a} r) \in \delta \wedge (w \rightarrow a) \in P)$ , 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  $\delta$  is acyclic, there exists a topological sort of  $\delta$  imposing a total order on  $Q$  such that  $M$  is strictly upper triangular (SUT). Initiate it thusly:

$$M_0[r, c, w] = \bigvee_{a \in \Sigma} \{a \mid (w \rightarrow a) \in P \wedge (q_r \xrightarrow{a} q_c) \in \delta\} \quad (2)$$

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

$$[\ell \oplus r]_w = [\ell_w \vee r_w] \quad (3)$$

$$[\ell \otimes r]_w = \bigvee_{x, z \in V} \{\ell_x \cdot r_z \mid (w \rightarrow xz) \in P\} \quad (4)$$

By slight abuse of notation<sup>2</sup>, we will redefine the matrix exponential over this domain as:

$$\exp(M) = \sum_{i=0}^{\infty} M_0^i = \sum_{i=0}^{|Q|} M_0^i \text{ (since } M \text{ is SUT.)} \quad (5)$$

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

$$T(2n) = \begin{cases} M_0, & \text{if } n = 1, \\ T(n) + T(n)^2 & \text{otherwise.} \end{cases} \quad (6)$$

Therefor, we only need at most  $\lceil \log_2 |Q| \rceil$  sequential steps to reach the fixpoint. Finally, we will union all the languages of every state pair deriving  $S$  into a new nonterminal,  $S_\cap$ .

$$S_\cap = \bigvee_{q \in I, q' \in F} \exp(M)[q, q', S] \text{ and } \varphi = [S_\cap \neq \emptyset] \quad (7)$$

To decode a witness in case of non-emptiness, one may simply choose  $(S_\cap)$ .  $\square$

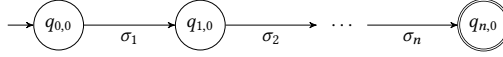
<sup>2</sup>Customarily, there is a  $\frac{1}{k!}$  factor to suppress exploding entries, but alas this domain has no multiplicative inverse.

## 6 EXAMPLES

In this section, we will consider three examples of intersections with finite languages. First, parsing can be viewed as a special case of language intersection with an automaton accepting a single word. Second, completion can be seen as a case of intersection with terminal wildcards in known locations. Thirdly, we consider syntax repair, where we will intersect a language representing all possible edit paths to determine the edit location(s) and fill them with appropriate terminals.

### 6.1 Recognition as intersection

In the case of ordinary CFL recognition, the automaton accepts just a single word:



Given a CFG,  $G' : \mathcal{G}$  in Chomsky Normal Form (CNF), we can construct a recognizer  $R : \mathcal{G} \rightarrow \Sigma^n \rightarrow \mathbb{B}$  for strings  $\sigma : \Sigma^n$  as follows. Let  $2^V$  be our domain,  $0$  be  $\emptyset$ ,  $\oplus$  be  $\cup$ , and  $\otimes$  be defined as:

$$X \otimes Z = \{ w \mid \langle x, z \rangle \in X \times Z, (w \rightarrow xz) \in P \} \quad (8)$$

If we define  $\hat{\sigma}_r = \{w \mid (w \rightarrow \sigma_r) \in P\}$ , then construct a matrix with nonterminals on the superdiagonal representing each token,  $M_0[r+1 = c](G', \sigma) = \hat{\sigma}_r$ , the fixpoint  $M_{i+1} = M_i + M_i^2$  is uniquely determined by the superdiagonal entries. Omitting the exponentiation-by-squaring detail, the ordinary fixedpoint iteration simply fills successive diagonals:

$$M_0 = \begin{pmatrix} \emptyset & \hat{\sigma}_1 & \emptyset & \dots & \emptyset \\ & \ddots & \ddots & \ddots & \ddots \\ \emptyset & \dots & \emptyset & \dots & \hat{\sigma}_n \\ & & & & \emptyset \end{pmatrix} \Rightarrow \begin{pmatrix} \emptyset & \hat{\sigma}_1 & \Lambda & \dots & \emptyset \\ & \ddots & \ddots & \ddots & \ddots \\ \emptyset & \dots & \emptyset & \dots & \hat{\sigma}_n \\ & & & & \emptyset \end{pmatrix} \Rightarrow \dots \Rightarrow M_\infty = \begin{pmatrix} \emptyset & \hat{\sigma}_1 & \Lambda & \dots & \Lambda_\sigma^* \\ & \ddots & \ddots & \ddots & \ddots \\ \emptyset & \dots & \emptyset & \dots & \hat{\sigma}_n \\ & & & & \emptyset \end{pmatrix}$$

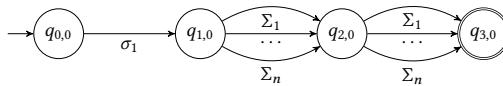
Once the fixpoint  $M_\infty$  is attained, the proposition  $[S \in \Lambda_\sigma^*]$  decides language membership, i.e.,  $[\sigma \in \mathcal{L}(G)]^3$ . So far, this procedure is essentially the textbook CYK algorithm in a linear algebraic notation [3] and a well-established technique in the parsing literature [4].

### 6.2 Completion as intersection

We can also consider a more general automaton for completing a string with holes, representing edits in fixed locations which can be filled by any terminal, which we call *completion*. In this case, the fixpoint is characterized by a system of language equations, whose solutions are the set of all sentences consistent with the template.

**Definition 6.1 (Completion).** Let  $\underline{\Sigma} = \Sigma \cup \{\_ \}$ , where  $\_$  denotes a hole. We denote  $\sqsubseteq : \Sigma^n \times \underline{\Sigma}^n$  as the relation  $\{\langle \sigma', \sigma \rangle \mid \sigma_i \in \Sigma \implies \sigma'_i = \sigma_i\}$  and the set of all inhabitants  $\{\sigma' : \Sigma^+ \mid \sigma' \sqsubseteq \sigma\}$  as  $H(\sigma)$ . Given a *porous string*,  $\sigma : \underline{\Sigma}^*$  we seek all syntactically valid inhabitants, i.e.,  $A(\sigma) = H(\sigma) \cap \ell$ .

Here, the FSA takes a similar shape but can have multiple arcs between subsequent states, e.g.:



<sup>3</sup>Hereinafter, we use Iverson brackets to denote the indicator function of a predicate with free variables, i.e.,  $[P] \Leftrightarrow \mathbb{1}(P)$ .

This corresponds to a template with two holes,  $\sigma = 1 \_ \_$ . Suppose the context-free grammar is  $G = \{S \rightarrow NON, O \rightarrow + \mid \times, N \rightarrow 0 \mid 1\}$ . This grammar will first be rewritten into CNF as  $G' = \{S \rightarrow NL, N \rightarrow 0 \mid 1, O \rightarrow \times \mid +, L \rightarrow ON\}$ . Using the powerset algebra we just defined, the matrix fixpoint  $M' = M + M^2$  can be computed as follows, shown in the leftmost column below:

	$2^V$	$\mathbb{Z}_2^{ V }$	$\mathbb{Z}_2^{ V } \rightarrow \mathbb{Z}_2^{ V }$
$M_0$	$\begin{pmatrix} \{N\} \\ \{N, O\} \\ \{N, O\} \end{pmatrix}$	$\begin{pmatrix} \begin{matrix} L & N & O & S \\ \square & \blacksquare & \square & \square \end{matrix} \\ \square & \blacksquare & \blacksquare & \square \\ \square & \blacksquare & \blacksquare & \square \end{pmatrix}$	$\begin{pmatrix} V_{0,1} \\ V_{1,2} \\ V_{2,3} \end{pmatrix}$
$M_1$	$\begin{pmatrix} \{N\} & \emptyset & \\ \{N, O\} & \{L\} & \\ & \{N, O\} & \end{pmatrix}$	$\begin{pmatrix} \square & \blacksquare & \square & \square & \square & \square & \square & \square \\ \square & \blacksquare & \blacksquare & \square & \blacksquare & \square & \square & \square \\ \square & \blacksquare & \blacksquare & \square & \blacksquare & \square & \square & \square \end{pmatrix}$	$\begin{pmatrix} V_{0,1} & V_{0,2} & \\ & V_{1,2} & V_{1,3} \\ & & V_{2,3} \end{pmatrix}$
$M_2$ = $M_\infty$	$\begin{pmatrix} \{N\} & \emptyset & \{S\} \\ \{N, O\} & \{L\} & \\ & \{N, O\} & \end{pmatrix}$	$\begin{pmatrix} \square & \blacksquare & \square & \square & \square & \square & \square & \blacksquare \\ \square & \blacksquare & \blacksquare & \square & \blacksquare & \square & \square & \square \\ \square & \blacksquare & \blacksquare & \square & \blacksquare & \square & \square & \square \end{pmatrix}$	$\begin{pmatrix} V_{0,1} & V_{0,2} & V_{0,3} \\ & V_{1,2} & V_{1,3} \\ & & V_{2,3} \end{pmatrix}$

The same procedure can be translated, without loss of generality, into the bit domain ( $\mathbb{Z}_2^{|V|}$ ) using a lexicographic nonterminal ordering, however  $M_\infty$  in both  $2^V$  and  $\mathbb{Z}_2^{|V|}$  represents a decision procedure, i.e.,  $[S \in V_{0,3}] \Leftrightarrow [V_{0,3,3} = \blacksquare] \Leftrightarrow [A(\sigma) \neq \emptyset]$ . Since  $V_{0,3} = \{S\}$ , we know there exists at least one solution  $\sigma' \in A(\sigma)$ , but  $M_\infty$  does not explicitly reveal its identity.

To extract the inhabitants, we can translate the bitwise procedure into an equation with free variables. Here, we can encode the idempotency constraint directly as  $M = M^2$ . We first define  $X \boxtimes Z = [X_2 \wedge Z_1, \perp, \perp, X_1 \wedge Z_0]$  and  $X \boxplus Z = [X_i \vee Z_i]_{i \in [0, |V|]}$ , mirroring  $\oplus, \otimes$  from the powerset domain, now over bitvectors. Since the unit nonterminals  $O, N$  can only occur on the superdiagonal, they may be safely ignored by  $\boxtimes$ . To solve for  $M_\infty$ , we proceed by first computing  $V_{0,2}, V_{1,3}$ :

$$\begin{aligned}
 V_{0,2} &= V_{0,j} \cdot V_{j,2} = V_{0,1} \boxtimes V_{1,2} & V_{1,3} &= V_{1,j} \cdot V_{j,3} = V_{1,2} \boxtimes V_{2,3} \\
 &= [L \in V_{0,2}, \perp, \perp, S \in V_{0,2}] & &= [L \in V_{1,3}, \perp, \perp, S \in V_{1,3}] \\
 &= [O \in V_{0,1} \wedge N \in V_{1,2}, \perp, \perp, N \in V_{0,1} \wedge L \in V_{1,2}] & &= [O \in V_{1,2} \wedge N \in V_{2,3}, \perp, \perp, N \in V_{1,2} \wedge L \in V_{2,3}] \\
 &= [V_{0,1,2} \wedge V_{1,2,1}, \perp, \perp, V_{0,1,1} \wedge V_{1,2,0}] & &= [V_{1,2,2} \wedge V_{2,3,1}, \perp, \perp, V_{1,2,1} \wedge V_{2,3,0}]
 \end{aligned}$$

Now we solve for the corner entry  $V_{0,3}$  by dotting the first row and last column, which yields:

$$\begin{aligned}
 V_{0,3} &= V_{0,j} \cdot V_{j,3} = (V_{0,1} \boxtimes V_{1,3}) \boxplus (V_{0,2} \boxtimes V_{2,3}) \\
 &= [V_{0,1,2} \wedge V_{1,3,1} \vee V_{0,2,2} \wedge V_{2,3,1}, \perp, \perp, V_{0,1,1} \wedge V_{1,3,0} \vee V_{0,2,1} \wedge V_{2,3,0}]
 \end{aligned}$$

Since we only care about  $V_{0,3,3} \Leftrightarrow [S \in V_{0,3}]$ , we can ignore the first three entries and solve for:

$$\begin{aligned}
 V_{0,3,3} &= V_{0,1,1} \wedge V_{1,3,0} \vee V_{0,2,1} \wedge V_{2,3,0} \\
 &= V_{0,1,1} \wedge (V_{1,2,2} \wedge V_{2,3,1}) \vee V_{0,2,1} \wedge \perp \\
 &= V_{0,1,1} \wedge V_{1,2,2} \wedge V_{2,3,1} \\
 &= [N \in V_{0,1}] \wedge [O \in V_{1,2}] \wedge [N \in V_{2,3}]
 \end{aligned}$$



### 6.3 Repair as intersection

Finally, we are ready to consider the general case of syntax repair, in which case the edit locations are not localized but can occur anywhere in the string. In this case, we construct a lattice of all possible edit paths within a fixed distance. This structure is called a Levenshtein automaton.

As the original construction defined by Schultz and Mihov [6] contains cycles and  $\varepsilon$ -transitions, we propose a variant which is  $\varepsilon$ -free and acyclic. Furthermore, we adopt a nominal form which supports infinite alphabets and considerably simplifies the language intersection to follow. Illustrated in Fig. 2 is an example of a small Levenshtein automaton recognizing  $L(\sigma : \Sigma^5, 3)$ . Unlabeled arcs accept any terminal from the alphabet,  $\Sigma$ . Equivalently, this transition system can be viewed as a kind of proof system within an unlabeled lattice. The following construction is equivalent to Schultz and Mihov's original Levenshtein automaton, but is more amenable to our purposes as it does not contain  $\varepsilon$ -arcs, and instead uses skip connections to recognize consecutive deletions of varying lengths.

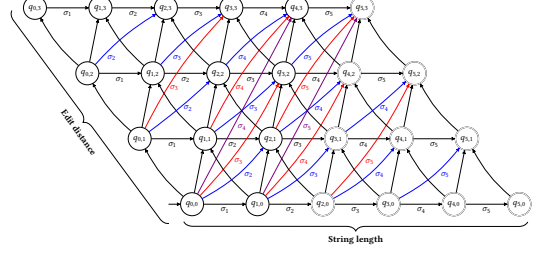
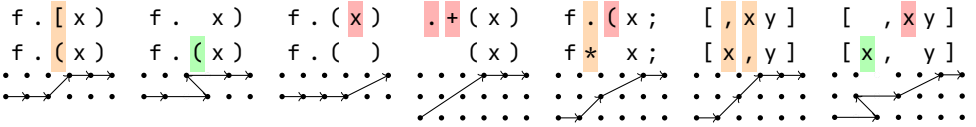


Fig. 2. NFA recognizing Levenshtein  $L(\sigma : \Sigma^5, 3)$ .

The following construction is equivalent to Schultz and Mihov's original Levenshtein automaton, but is more amenable to our purposes as it does not contain  $\varepsilon$ -arcs, and instead uses skip connections to recognize consecutive deletions of varying lengths.

$$\begin{array}{c}
 \frac{s \in \Sigma \quad i \in [0, n] \quad j \in [1, d_{\max}]}{(q_{i,j-1} \xrightarrow{s} q_{i,j}) \in \delta} \quad \nwarrow \quad \frac{s \in \Sigma \quad i \in [1, n] \quad j \in [1, d_{\max}]}{(q_{i-1,j-1} \xrightarrow{s} q_{i,j}) \in \delta} \quad \nearrow \\
 \frac{i \in [1, n] \quad j \in [0, d_{\max}]}{(q_{i-1,j} \xrightarrow{\sigma_i} q_{i,j}) \in \delta} \quad \rightarrow \quad \frac{d \in [1, d_{\max}] \quad i \in [d+1, n] \quad j \in [d, d_{\max}]}{(q_{i-d-1,j-d} \xrightarrow{\sigma_i} q_{i,j}) \in \delta} \quad \nearrow \\
 \frac{}{q_{0,0} \in I} \text{ INIT} \quad \frac{q_{i,j} \in Q \quad |n-i+j| \leq d_{\max}}{q_{i,j} \in F} \text{ DONE}
 \end{array}$$

Each arc plays a specific role.  $\nwarrow$  handles insertions,  $\nearrow$  handles substitutions and  $\nearrow$  handles deletions of one or more terminals. Let us consider some illustrative cases.



Note that the same patch can have multiple Levenshtein alignments. DONE constructs the final states, which are all states accepting strings  $\sigma'$  whose Levenshtein distance  $\Delta(\sigma, \sigma') \leq d_{\max}$ .

To avoid creating a parallel bundle of arcs for each insertion and substitution point, we instead decorate each arc with a nominal predicate, accepting or rejecting  $\sigma_i$ . To distinguish this nominal variant from the original construction, we highlight the modified rules in orange below.

$$\begin{array}{c}
 \frac{i \in [0, n] \quad j \in [1, d_{\max}]}{(q_{i,j-1} \xrightarrow{[\neq \sigma_{i+1}]} q_{i,j}) \in \delta} \quad \nwarrow \quad \frac{i \in [1, n] \quad j \in [1, d_{\max}]}{(q_{i-1,j-1} \xrightarrow{[\neq \sigma_i]} q_{i,j}) \in \delta} \quad \nearrow \\
 \frac{i \in [1, n] \quad j \in [0, d_{\max}]}{(q_{i-1,j} \xrightarrow{[= \sigma_i]} q_{i,j}) \in \delta} \quad \rightarrow \quad \frac{d \in [1, d_{\max}] \quad i \in [d+1, n] \quad j \in [d, d_{\max}]}{(q_{i-d-1,j-d} \xrightarrow{[\neq \sigma_i]} q_{i,j}) \in \delta} \quad \nearrow
 \end{array}$$

Nominalizing the NFA eliminates the creation of  $e = 2(|\Sigma| - 1) \cdot |\sigma| \cdot d_{\max}$  unnecessary arcs over the entire Levenshtein automaton and drastically reduces the representation size, but does not affect the underlying semantics. Thus, it is important to first nominalize the automaton before proceeding to avoid a large blowup in the intermediate grammar.

As a concrete example, suppose we have the string,  $\sigma = ( ) )$  and wish to balance the parentheses. We will initially have the Levenshtein automaton,  $A$ , depicted in Fig. 3. To check for non-emptiness, we will perform the following procedure. Suppose we have a CNF CFG,  $G' = \{S \rightarrow LR, S \rightarrow LF, S \rightarrow SS, F \rightarrow SR, L \rightarrow (, R \rightarrow )\}$  and let us assume an ordering of  $S, F, L, R$  on  $V$ .

First, we need to order the automata states by increasing longest-path distance from  $q_0$ . One approach would be to topologically sort the adjacency matrix. While some form of sorting is unavoidable for arbitrary ANFAs, if we know ahead of time that our structure is a Levenshtein automaton, we can simply enumerate its state space by increasing Manhattan distance from the origin. So, a valid ordering on  $Q$  would be  $q_{00}, q_{01}, q_{10}, q_{11}, q_{20}, q_{21}, q_{30}, q_{31}$ . Now, we want to compute  $L(G') \cap L(A)$ .

Under such an ordering, the adjacency matrix takes an upper triangular form. This forms the template for the initial parse chart. The initial parse chart,  $M_0$  (Fig. 6) and its fixpoint  $M_\infty$ , post convergence, (Fig. 7) are shown. Each entry of this chart corresponds to a vector of expressions  $E^{|V|}$  with at least one expression denoting a nonempty language. The adjacency and reachability matrices will always cover the expression vectors of the initial and final parse charts. In other words, we can always skip empty states in the reachability matrix.

Since  $M_\infty[q_{00}, q_{31}, 0] = \blacksquare$ , this implies that  $L(A) \cap L(G') \neq \emptyset$ , hence  $\text{LED}(\sigma, G) = 1$ . Using the same matrix, we can then perform a second pass over the nonempty entries to construct regular expression vectors representing finite languages for each constituent and decode a witness using the Brzowski derivative.

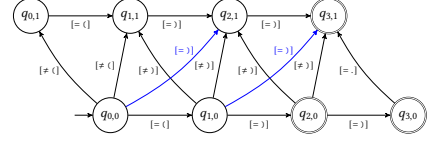


Fig. 3. Simple Levenshtein automaton.

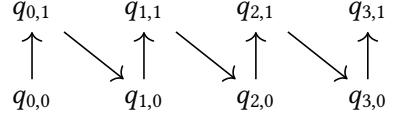


Fig. 4. Pairing function over  $L(\sigma : \Sigma^3, 1)$ .

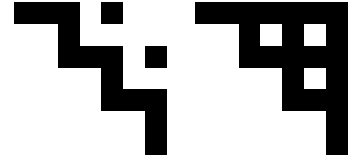


Fig. 5. Adjacency and reachability matrix.

$M_0$	$q_{00}$	$q_{01}$	$q_{10}$	$q_{11}$	$q_{20}$	$q_{21}$	$q_{30}$	$q_{31}$
$q_{00}$	S F L R	S F L R	S F L R	S F L R	S F L R	S F L R	S F L R	S F L R
$q_{01}$								
$q_{10}$								
$q_{11}$								
$q_{20}$								
$q_{21}$								
$q_{30}$								
$q_{31}$								

Fig. 6. Initial parse chart configuration.

$M_\infty$	$q_{00}$	$q_{01}$	$q_{10}$	$q_{11}$	$q_{20}$	$q_{21}$	$q_{30}$	$q_{31}$
$q_{00}$	S F L R	S F L R	S F L R	S F L R	S F L R	S F L R	S F L R	S F L R
$q_{01}$								
$q_{10}$								
$q_{11}$								
$q_{20}$								
$q_{21}$								
$q_{30}$								
$q_{31}$								

Fig. 7. Final parse chart configuration.

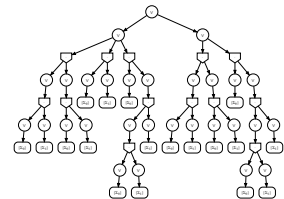


Fig. 8. Regular expression denoting  $\mathcal{L}(G_N)$ .

## 7 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.

### 7.1 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.

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

$$\text{Definition 7.1 (Cardinality). } |e| : E \rightarrow \mathbb{N} = \begin{cases} 1 & \text{if } R \in \Sigma \\ x \times z & \text{if } e = x \cdot z \\ x + z & \text{if } e = x \vee z \end{cases}$$

THEOREM 7.2 (ENUMERATION). To enumerate, invoke  $\bigcup_{i=0}^{|R|} \{enum(R, i)\}$ :

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

### 7.2 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:

$$S \rightarrow A B \ (0.9999) \quad S \rightarrow C C \ (0.0001)$$

$$A \rightarrow a \ (1) \quad B \rightarrow b \ (1) \quad C \rightarrow a \left(\frac{1}{26}\right) \mid \dots \mid z \left(\frac{1}{26}\right)$$

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.

### 7.3 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 7.3. If the FSA,  $\alpha$ , is ambiguous, then the intersection grammar,  $G_{\cap}$ , can be ambiguous.

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

REFERENCES

[1] Janusz A Brzozowski. 1964. Derivatives of regular expressions. Journal of the ACM (JACM) 11, 4 (1964), 481–494.

[2] Noam Chomsky. 1959. On certain formal properties of grammars. Information and control 2, 2 (1959), 137–167.

[3] Joshua Goodman. 1999. Semiring parsing. Computational Linguistics 25, 4 (1999), 573–606. <https://aclanthology.org/J99-4004.pdf>

[4] Dick Grune and Ceriel J. H. Jacobs. 2008. Parsing as Intersection. Springer New York, New York, NY, 425–442. [https://doi.org/10.1007/978-0-387-68954-8\\_13](https://doi.org/10.1007/978-0-387-68954-8_13)

[5] Arto Salomaa. 1973. Formal languages. Academic Press, New York. 59–61 pages.

[6] Klaus U Schulz and Stoyan Mihov. 2002. Fast string correction with Levenshtein automata. International Journal on Document Analysis and Recognition 5 (2002), 67–85.

## A LEVENSHTTEIN AUTOMATA MATRICES

These are useful for visually checking different implementations.

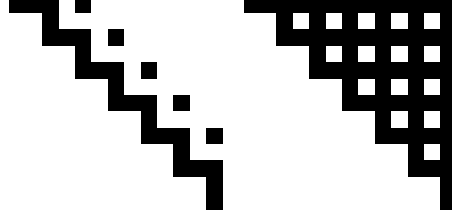


Fig. 9.  $\text{Lev}(|\sigma|=6, \Delta=1)$  adjacency and reachability matrices.



Fig. 10.  $\text{Lev}(|\sigma|=6, \Delta=2)$  adjacency and reachability matrices.

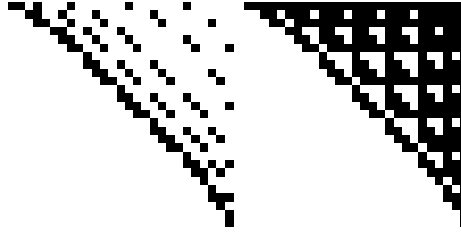


Fig. 11.  $\text{Lev}(|\sigma|=6, \Delta=3)$  adjacency and reachability matrices.

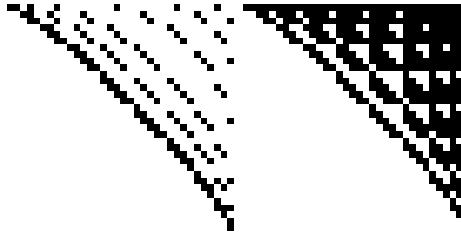


Fig. 12.  $\text{Lev}(|\sigma|=6, \Delta=4)$  adjacency and reachability matrices.

## B LEVENSHTEIN AUTOMATA MINIMALITY

It is reasonable to ask whether the Levenshtein automaton defined is minimal, in the sense of whether there exists an automaton with fewer states than  $A$  yet still generates  $\mathcal{L}(G_\cap)$  when intersected with  $\mathcal{L}(G)$ . In other words, given  $G$  and  $\sigma$ , is there an  $A'$  such that  $|Q_{A'}| < |Q_A|$  yet  $\mathcal{L}(G) \cap \mathcal{L}(A') = \mathcal{L}(G) \cap \mathcal{L}(A)$  still holds? In fact, there is a trivial example:

**THEOREM B.1.** *Let  $Q_{A'}$  be defined as  $Q_A \setminus \{q_{n,0}\}$ .*

Since  $q_{n,0}$  accepts the original string  $\sigma$  which by definition is not in  $\mathcal{L}(G)$ , we can immediately rule out this state. Moreover, we can define a family of automata with strictly fewer states than the full LBH construction by making the following observation: if we can prove one edit must occur before the last  $s$  tokens, we can rule out the last  $s$  states absorbing editless trajectories.

**THEOREM B.2.**  *$\emptyset = \mathcal{L}(\sigma_{1\dots(n-s)} \cdot \Sigma^s) \cap \mathcal{L}(G)$  implies the states  $[q_{n-i,0}]_{i=1\dots s}$  are unnecessary.*

Likewise, if we expend our entire edit budget in the first  $p$  tokens, we will be unable to recover in a string where at least one repair must occur after the first  $p$  tokens.

**THEOREM B.3.**  *$\emptyset = \mathcal{L}(\Sigma^p \cdot \sigma_p) \cap \mathcal{L}(G)$  implies the states  $[q_{i,d_{\max}}]_{i=0\dots p}$  are unnecessary.*

Therefor, we can eliminate  $p+s$  states from  $A$  by proving emptiness of  $\mathcal{L}(\Sigma^p \cdot \sigma_{p\dots(n-s)} \cdot \Sigma^s) \cap \mathcal{L}(G)$ , without affecting  $\mathcal{L}(G_\cap)$ . Pictorially,

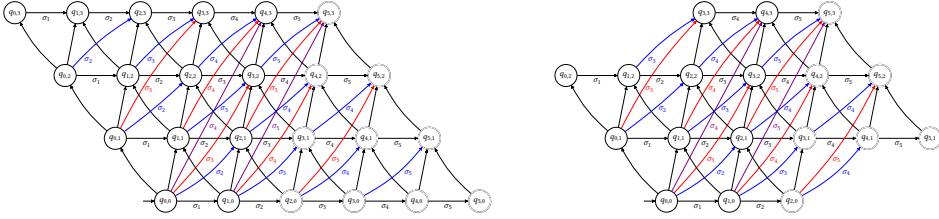


Fig. 13. Levenshtein NFA before and after pruning.

Pruned L-NFA for the broken string  $\sigma = [ ( + ) ]$  with  $G = \{S \rightarrow (S) \mid [S] \mid S + S \mid 1\}$ .

$- \quad - \quad + \quad ) \quad ]$	$\times$	$\wedge$	$- \quad - \quad - \quad ) \quad ]$	$\checkmark$
$[ \quad ( \quad + \quad - \quad - \quad ]$	$\times$	$\wedge$	$[ \quad ( \quad - \quad - \quad - \quad ]$	$\checkmark$