

# Syntax Repair as Idempotent Tensor Completion

Breandan Considine<sup>1</sup>, Jin Guo<sup>1</sup>, and Xujie Si<sup>2</sup>

<sup>1</sup> McGill University, Montréal, QC H2R 2Z4, Canada  
{breandan.considine@mail, jguo@cs}.mcgill.ca

<sup>2</sup> University of Toronto, Toronto, ON, M5S 1A1 Canada  
six@utoronto.ca

**Abstract.** We introduce a new technique for correcting syntax errors in arbitrary context-free languages. To do so, we reduce CFL recognition onto a Boolean tensor completion and compare various techniques for introducing the holes, and solving for their inhabitants. Our technique has practical applications for real-time syntax error correction.

**Keywords:** Error correction · CFL reachability · Language games.

## 1 Introduction

Syntax repair is the problem of taking a grammar and a malformed string, and modifying the string so it conforms to the grammar. Prior work has been devoted to fixing syntax errors using handcrafted heuristics. We take a first-principles approach that makes no assumptions about the string or grammar and focuses on accuracy and end-to-end latency. The result is a tool that is applicable to any context-free and conjunctive language, and which is provably sound and complete up to a Levenshtein bound.

### 1.1 Problem

Syntax repair can be treated as a language intersection problem between a context-free language (CFL) and a regular language.

**Definition 1 (Bounded Levenshtein-CFL reachability).** *Given a CFL  $\ell$  and an invalid string  $\underline{\sigma} : \ell^{\mathbb{C}}$ , the BCFLR problem is to find every valid string reachable within  $d$  edits of  $\underline{\sigma}$ , i.e., letting  $\Delta$  be the Levenshtein metric and  $L(\underline{\sigma}, d) := \{\sigma \mid \Delta(\underline{\sigma}, \sigma) \leq d\}$ , we seek to find  $L(\underline{\sigma}, d) \cap \ell$ .*

To solve this problem, we will first pose a simpler problem that only considers intersections with a finite language, then turn our attention back to BCFLR.

**Definition 2 (Porous completion).** *Let  $\underline{\Sigma} := \Sigma \cup \{\_\}$ , where  $\_$  denotes a hole. We denote  $\sqsubseteq : \Sigma^n \times \underline{\Sigma}^n$  as the relation  $\{(\sigma', \sigma) \mid \sigma_i \in \Sigma \implies \sigma'_i = \sigma_i\}$  and the set of all inhabitants  $\{\sigma' \mid \sigma' \sqsubseteq \sigma\}$  as  $H(\sigma)$ . Given a porous string,  $\sigma : \underline{\Sigma}^*$  we seek all syntactically admissible inhabitants, i.e.,  $A(\sigma) := H(\sigma) \cap \ell$ .*

$A(\sigma)$  is often a large-cardinality set, so we want a procedure which returns the most likely members first, without exhaustive enumeration. More precisely,

**Definition 3 (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- $k$  repairs by likelihood under the language model. That is,*

$$R(\ell^\cap, P_\theta) := \underset{\{\sigma \mid \sigma \subseteq \ell^\cap, |\sigma| \leq k\}}{\operatorname{argmax}} \sum_{\sigma \in \sigma} \prod_{i=1}^{|\sigma|} P_\theta(\sigma_i \mid \sigma_{1\dots i})^{\frac{1}{|\sigma|}} \quad (1)$$

We want a procedure  $\hat{R}$ , minimizing  $\mathbb{E}_{G, \sigma} [D_{KL}(\hat{R} \parallel R)]$  and wallclock runtime.

Our key innovation and the core problem this paper tackles is, given  $\underline{\sigma}, d, P_\theta$ , to approximate  $R(\ell^\cap, P_\theta)$  while minimizing latency and maximizing accuracy. We will start with some background, give an example, then dive into the theory.

## 1.2 Background

Recall that a CFG 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 *Chomsky Normal Form*,  $P': V \rightarrow (V^2 \mid \Sigma)$ , in which every  $P$  takes one of two forms, either  $w \rightarrow xz$ , or  $w \rightarrow t$ , where  $w, x, z : V$  and  $t : \Sigma$ . For example:

$$G := \{ S \rightarrow S S \mid (S) \mid () \} \implies \{ S \rightarrow Q R \mid S S \mid L R, R \rightarrow ), L \rightarrow (, Q \rightarrow L S \}$$

Given a CFG,  $G' : \mathbb{G} = \langle \Sigma, V, P, S \rangle$  in CNF, we can construct a recognizer  $R : \mathbb{G} \rightarrow \Sigma^n \rightarrow \mathbb{B}$  for strings  $\sigma : \Sigma^n$  as follows. Let  $2^V$  be our domain, 0 be  $\emptyset$ ,  $\oplus$  be  $\cup$ , and  $\otimes$  be defined as:

$$X \otimes Z := \{ w \mid \langle x, z \rangle \in X \times Z, (w \rightarrow xz) \in P \} \quad (2)$$

If we define  $\sigma_r := \{ w \mid (w \rightarrow \sigma_r) \in P \}$ , then construct a matrix with nonterminals on the superdiagonal representing each token,  $M_{r+1=c}(G', e) := \sigma_r$  and solve for the fixpoint  $M_{i+1} = M_i + M_i^2$ ,

$$M_0 := \begin{pmatrix} \emptyset & \sigma_1 & \emptyset & \emptyset \\ & \ddots & \ddots & \emptyset \\ & & \emptyset & \sigma_n \\ \emptyset & \dots & \dots & \emptyset \end{pmatrix} \Rightarrow \begin{pmatrix} \emptyset & \sigma_1 & \Lambda & \emptyset \\ & \ddots & \ddots & \emptyset \\ & & \Lambda & \sigma_n \\ \emptyset & \dots & \dots & \emptyset \end{pmatrix} \Rightarrow \dots \Rightarrow M_\infty = \begin{pmatrix} \emptyset & \sigma_1 & \Lambda & \Lambda_\sigma^* \\ & \ddots & \ddots & \emptyset \\ & & \Lambda & \sigma_n \\ \emptyset & \dots & \dots & \emptyset \end{pmatrix}$$

we obtain the recognizer,  $R(G', \sigma) := [S \in \Lambda_\sigma^*] \Leftrightarrow [\sigma \in \mathcal{L}(G)]$ <sup>3</sup>.

Since  $\bigoplus_{c=1}^n M_{r,c} \otimes M_{c,r}$  has cardinality bounded by  $|V|$ , it can be represented as  $\mathbb{Z}_2^{|V|}$  using the characteristic function,  $\mathbb{1}$ . A concrete example is shown in § 1.3.

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

### 1.3 Example

Let us consider an example with two holes,  $\sigma = 1 \_ \_$ , and the grammar being  $G := \{S \rightarrow NON, O \rightarrow + \mid \times, N \rightarrow 0 \mid 1\}$ . This can be rewritten into CNF as  $G' := \{S \rightarrow NL, N \rightarrow 0 \mid 1, O \rightarrow \times \mid +, L \rightarrow ON\}$ . Using the algebra where  $\oplus := \cup$ ,  $X \otimes Z := \{ w \mid \langle x, z \rangle \in X \times Z, (w \rightarrow xz) \in P \}$ , the fixpoint  $M' = M + M^2$  can be computed as follows, shown in the leftmost column:

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

The same procedure can be translated, without loss of generality, into the bit domain  $(\mathbb{Z}_2^{|V|})$  using a lexicographic ordering, however these both are recognizers. That is to say,  $[S \in V_{0,3}] \Leftrightarrow [V_{0,3,3} = 1] \Leftrightarrow [A(\sigma) \neq \emptyset]$ . Since  $V_{0,3} = \{S\}$ , we know there is at least one  $\sigma' \in A(\sigma)$ , but  $M_\infty$  does not reveal its identity.

In order 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  $X \boxplus Z := [X_i \vee Z_i]_{i \in [0, |V|]}$ . Since the unit nonterminals  $O, N$  can only occur on the superdiagonal, they may be safely ignored by  $\otimes$ . To solve for  $M_\infty$ , we proceed by first computing  $V_{0,2}, V_{0,3}$  as follows:

$$V_{0,2} = V_{0,j} \cdot V_{j,2} = V_{0,1} \boxtimes V_{1,2} \quad (3)$$

$$= [L \in V_{0,2}, \perp, \perp, S \in V_{0,2}] \quad (4)$$

$$= [O \in V_{0,1} \wedge N \in V_{1,2}, \perp, \perp, N \in V_{0,1} \wedge L \in V_{1,2}] \quad (5)$$

$$= [V_{0,1,2} \wedge V_{1,2,1}, \perp, \perp, V_{0,1,1} \wedge V_{1,2,0}] \quad (6)$$

$$V_{1,3} = V_{1,j} \cdot V_{j,3} = V_{1,2} \boxtimes V_{2,3} \quad (7)$$

$$= [L \in V_{1,3}, \perp, \perp, S \in V_{1,3}] \quad (8)$$

$$= [O \in V_{1,2} \wedge N \in V_{2,3}, \perp, \perp, N \in V_{1,2} \wedge L \in V_{2,3}] \quad (9)$$

$$= [V_{1,2,2} \wedge V_{2,3,1}, \perp, \perp, V_{1,2,1} \wedge V_{2,3,0}] \quad (10)$$

Now we solve for the corner entry  $V_{0,3}$  by taking the bitwise dot product between the first row and last column, yielding:

$$V_{0,3} = V_{0,j} \cdot V_{j,3} = V_{0,1} \boxtimes V_{1,3} \boxplus V_{0,2} \boxtimes V_{2,3} \quad (11)$$

$$= [V_{0,1,2} \wedge V_{1,3,1} \vee V_{0,2,2} \wedge V_{2,3,1}, \perp, \perp, V_{0,1,1} \wedge V_{1,3,0} \vee V_{0,2,1} \wedge V_{2,3,0}] \quad (12)$$

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

$$V_{0,3,3} = V_{0,1,1} \wedge V_{1,3,0} \vee V_{0,2,1} \wedge V_{2,3,0} \quad (13)$$

$$= V_{0,1,1} \wedge (V_{1,2,2} \wedge V_{2,3,1}) \vee V_{0,2,1} \wedge \perp \quad (14)$$

$$= V_{0,1,1} \wedge V_{1,2,2} \wedge V_{2,3,1} \quad (15)$$

$$= [N \in V_{0,1}] \wedge [O \in V_{1,2}] \wedge [N \in V_{2,3}] \quad (16)$$

Now we know that  $\sigma = 1 \ \underline{O} \ \underline{N}$  is a valid solution, and therefore we can take the product  $\{1\} \times \sigma_r^{-1}(O) \times \sigma_r^{-1}(N)$  to recover the admissible set, yielding  $A(\sigma) = \{1 + 0, 1 + 1, 1 \times 0, 1 \times 1\}$ . In this case, since  $G$  is unambiguous, there is only one parse tree satisfying  $V_{0,|\sigma|,|\sigma|}$ , but in general, there can be multiple valid parse trees, in which case we can decode them incrementally.

#### 1.4 Semiring Algebras

There are a number of alternate semirings which solve the same problem. A first approach requires solving for  $A(\sigma)$  using a semiring algebra, and propagating the values from the bottom-up as a string to a list of strings. Letting  $D = V \rightarrow \mathcal{P}(\Sigma^*)$ , we define  $\oplus, \otimes : D \times D \rightarrow D$ . Initially, we have  $p(s : \Sigma) := \{v \mid [v \rightarrow s] \in P\}$  and  $p(\_) := \bigcup_{s \in \Sigma} p(s)$ , then we compute the fixpoint using the algebra:

$$X \oplus Z := \{v \rightarrow (X(v) \cup Z(v)) \mid v \in V\} \quad (17)$$

$$X \otimes Z := \bigoplus_{w,x,z} \{w \rightarrow (l+r) \mid [w \rightarrow xz] \in P, \langle l, r \rangle \in X(x) \times Z(x)\} \quad (18)$$

After the fixpoint  $M_\infty$  is attained, the solutions can be read off via  $M_\infty[0, |\sigma|](S)$ . The issue here is an exponential growth in cardinality when eagerly computing the Cartesian product, which becomes impractical for even small strings. We can make this encoding more compact by propagating an algebraic data type (ADT)  $\mathbb{T}_2$  using the operations  $\oplus, \otimes : 2^{\mathbb{T}_2} \times 2^{\mathbb{T}_2} \rightarrow 2^{\mathbb{T}_2}$  as follows:

$$X \oplus Z := \{\mathbb{T}_2(k, Q_x \cup Q_z) \mid (k, Q)_X \bowtie_k (k, Q)_Z\} \quad (19)$$

$$X \otimes Z := \bigoplus_{w,x,z} \{\mathbb{T}_2(w, \{\langle T_x, T_z \rangle\}) \mid [w \rightarrow xz] \in P, x \in \pi_1(X), z \in \pi_1(Z)\} \quad (20)$$

Decoding, then becomes a matter of enumerating binary trees from the ADT using a recursive choice function that emits a sequence of strings satisfying  $A(\sigma)$ , with the type signature  $\mathcal{C} : \mathbb{T}_2 \rightarrow (\mathbb{N} \rightarrow \Sigma^*)$  defined as follows:

$$\mathcal{C}(t : \mathbb{T}_2) := \begin{cases} \pi_1(t) & \text{if } \pi_2(t) = \emptyset, \text{ or} \\ \{x + z \mid \langle X, Z \rangle \in \pi_2(t), x \in \mathcal{C}(X), z \in \mathcal{C}(Z)\} & \text{otherwise.} \end{cases}$$

In our experiments, we provide a comparison of the performance of the SAT algebra and these two semirings, evaluated on a dataset of Python statements.

### 1.5 Complexity

Let us consider some loose bounds on the complexity of BCFL. To do, we first consider the complexity of computing language-edit distance, which is a lower-bound on the complexity of BCFL.

**Definition 4.** *Language edit distance (LED) is the problem of computing the minimum number of edits required to transform an invalid string into a valid one, where validity is defined as containment in a context-free language,  $\ell : \mathcal{L}$ , i.e.,  $\Delta^*(\sigma, \ell) := \min_{\sigma \in \ell} \Delta(\sigma, \sigma)$ , and  $\Delta$  is the Levenshtein distance. LED is known to have a subcubic complexity. [1]*

We seek to find the set of strings  $S$  such that  $\forall \tilde{\sigma} \in S, \Delta(\sigma, \tilde{\sigma}) \leq q$ , where  $q$  is the maximum number of edits greater than or equal to the language edit distance. We call this set the *Levenshtein ball* of  $\sigma$  and denote it  $\Delta_q(\sigma)$ . Since  $1 \leq \Delta^*(\sigma, \ell) \leq q$ , we have  $1 \leq q$ . We now consider an upper bound on  $\Delta^*(\sigma, \ell)$ , i.e., the greatest lower bound on  $q$  such that  $\Delta_q(\sigma) \cap \ell \neq \emptyset$ .

**Lemma 1.** *For any nonempty language  $\ell : \mathcal{L}$  and invalid string  $\sigma : \Sigma^n \cap \bar{\ell}$ , there exists an  $(\tilde{\sigma}, m)$  such that  $\tilde{\sigma} \in \ell \cap \Sigma^m$  and  $0 < \Delta(\sigma, \ell) \leq \max(m, n) < \infty$ .*

*Proof.* Since  $\ell$  is nonempty, it must have at least one inhabitant  $\sigma \in \ell$ . Let  $\tilde{\sigma}$  be the smallest such member. Since  $\tilde{\sigma}$  is a valid sentence in  $\ell$ , by definition it must be that  $|\tilde{\sigma}| < \infty$ . Let  $m := |\tilde{\sigma}|$ . Since we know  $\sigma \notin \ell$ , it follows that  $0 < \Delta(\sigma, \ell)$ . Let us consider two cases, either  $\tilde{\sigma} = \varepsilon$ , or  $0 < |\tilde{\sigma}|$ :

- If  $\tilde{\sigma} = \varepsilon$ , then  $\Delta(\sigma, \tilde{\sigma}) = n$  by full erasure of  $\sigma$ , or
- If  $0 < m$ , then  $\Delta(\sigma, \tilde{\sigma}) \leq \max(m, n)$  by overwriting.

In either case, it follows  $\Delta(\sigma, \ell) \leq \max(m, n)$  and  $\ell$  is always reachable via a finite nonempty set of Levenshtein edits, i.e.,  $0 < \Delta(\sigma, \ell) < \infty$ .

Let us now consider the maximum growth rate of the *admissible set*,  $A := \Delta_q(\sigma) \cap \ell$ , as a function of  $q$  and  $n$ . Let  $\bar{\ell} := \{\sigma\}$ . Since  $\bar{\ell}$  is finite and thus regular,  $\ell = \Sigma^* \setminus \{\sigma\}$  is regular by the closure of regular languages under complementation, and thus context-free a fortiori. Since  $\ell$  accepts every string except  $\sigma$ , it represents the worst CFL in terms of asymptotic growth of  $A$ .

**Lemma 2.** *BCFL complexity is upper bounded by  $\mathcal{O}(\sum_{c=1}^q \binom{cn+n+1}{c} (|\Sigma| + 1)^c)$ .*

*Proof.* We can overestimate the size of  $A$  by considering the number of unique ways to insert, delete, or substitute  $c$  terminals into a string  $\underline{\sigma}$  of length  $n$ . This can be overapproximated by interleaving  $\varepsilon^c$  around every token, i.e.,  $\underline{\sigma}_\varepsilon := (\varepsilon^c \underline{\sigma}_i)_{i=1}^n \varepsilon^c$ , where  $|\underline{\sigma}_\varepsilon| = cn + n + 1$ , and only considering substitution. We augment  $\Sigma_\varepsilon := \Sigma \cup \{\varepsilon\}$  so that deletions and insertions may be treated as special cases of substitution. Thus, we have  $cn + n + 1$  positions to substitute  $(|\Sigma_\varepsilon|)$  tokens, i.e.,  $\binom{cn+n+1}{c} |\Sigma_\varepsilon|^c$  ways to edit  $\underline{\sigma}_\varepsilon$  for each  $c \in [1, q]$ . This upper bound is not tight, as overcounts many identical edits w.r.t.  $\underline{\sigma}$ . Nonetheless, it is sufficient to show  $|A| < \sum_{c=1}^q \binom{cn+n+1}{c} |\Sigma_\varepsilon|^c$ .

We note that the above bound applies to all strings and languages, and relates to the Hamming bound on  $H_q(\underline{\sigma}_\varepsilon)$ , which only considers substitutions.<sup>4</sup> In practice, much tighter bounds may be obtained by considering the structure of  $\ell$  and  $\underline{\sigma}$ . For example, based on an empirical evaluation from a dataset of human errors and repairs in Python code snippets ( $|\Sigma| = 50, |\underline{\sigma}| < 40, \Delta(\underline{\sigma}, \ell) \in [1, 3]$ ), we estimate the *filtration rate*, i.e., the density of the admissible set relative to the Levenshtein ball,  $D = |A|/|\Delta_q(\underline{\sigma})|$  to have empirical mean  $E_\sigma[D] \approx 2.6 \times 10^{-4}$ , and variance  $\text{Var}_\sigma[D] \approx 3.8 \times 10^{-7}$ .

In practice, this problem is ill-posed even when  $q = \Delta^*(\underline{\sigma}, \ell) \approx 1$ . For example, consider the language of ursine dietary preferences. Although  $\underline{\sigma} :=$  “Bears like to eat plastic” is not a valid sentence, e.g.,  $\tilde{\sigma} :=$  “Bears like to eat” is ( $\Delta^* = 1$ ), however there are many others with roughly the same edit distance, e.g., “Bears like to eat {berries, honey, fish}”, or “{Polar, Panda} bears like to eat {seals, bamboo}”. In general, there are usually many strings nearby  $\underline{\sigma}$ , and we seek to find those among them which are both syntactically valid and semantically plausible as quickly as possible.

## 1.6 Our solution to bounded CFL reachability

Now, let us return to the problem of bounded CFL reachability. The question naturally arises, of should one put the holes? One solution is to sample holes without replacement from all possible locations. Below is an example illustrating this procedure on a single Python snippet.

```

1. d = sum([foo(i) for i in vals])
2. d = sum([ [foo([i] for i in vals)] ])
3. w = w([ [w([w] for w in w)] ])
4. [w] = [w] ([ [w] ([w] ] for [w] in [w] ) ) ]
5. [ ] = [ ] ([ [w] ([w] ] for [w] in [w] ) ) [ ]
   [w] [ ] [w] ([ [ ] ([w] ] for [w] [ ] [w] ) ) ]
...

```

<sup>4</sup> This reflects our general approach, which builds a surjection from the interleaved Hamming ball onto the Levenshtein ball.

6. 

w	=	w	(	[	w	(	-	w	-	for	w	in	w	-	)
							+	)						]	

  
7. `d = sum([foo(+i) for i in vals])`
  
8. `d = sum([foo(i) for i in vals])`

The initial broken string, `d = sum([foo(i) for i in vals])` (1), is first tokenized using a lexer to obtain the sequence in (2). Lexical tokens containing identifiers are abstracted in step (3), and interleaved with the empty token in step (4). We then sample hole configurations without replacement in step (5), many of which will have no admissible solutions. Eventually, the solver will discover an admissible solution, as seen in step (6). This solution is then used to generate a patch, which is applied to the original string in step (7), then reduced to its minimal form in step (8), and sampling is repeated until all possibilities are exhausted or a predetermined timeout expires.

### 1.7 Sampling the Levenshtein ball without replacement in $\mathcal{O}(1)$

Now that we have a reliable method to synthesize admissible completions for strings containing holes, i.e., fix *localized* errors,  $F : (\mathcal{G} \times \Sigma^n) \rightarrow \{\Sigma^n\} \subseteq \mathcal{L}(\mathcal{G})$ , how can we use  $F$  to repair some unparseable string, i.e.,  $\sigma_1 \dots \sigma_n : \Sigma^n \cap \mathcal{L}(\mathcal{G})^c$  where the holes' locations are unknown? Three questions stand out in particular: how many holes are needed to repair the string, where should we put those holes, and how ought we fill them to obtain a parseable  $\tilde{\sigma} \in \mathcal{L}(\mathcal{G})$ ?

One plausible approach would be to draw samples with a PCFG, minimizing tree-edit distance, however these are computationally expensive metrics and approximations may converge poorly. A more efficient strategy is to sample string perturbations,  $\sigma \sim \Sigma^{n \pm q} \cap \Delta_q(\underline{\sigma})$  uniformly across the Levenshtein  $q$ -ball centered on  $\underline{\sigma}$ , i.e., the space of all admissible edits with Levenshtein distance  $\leq q$ .

To implement this strategy, we first construct a surjection  $\varphi^{-1} : \mathbb{Z}_2^m \twoheadrightarrow \Delta_q(\underline{\sigma})$  from bitvectors to Levenshtein edits over  $\underline{\sigma}$ ,  $\Sigma$ , sample bitvectors without replacement using a characteristic polynomial, then decode the resulting bitvectors into Levenshtein edits. This ensures the sampler eventually visits every Levenshtein edit at least exactly once and at most approximately once, without needing to store any samples, and discovers a steady stream of admissible edits throughout the solving process, independent of the grammar or string under repair.

More specifically, we employ a pair of [un]tupling functions  $\kappa, \rho : \mathbb{N}^k \leftrightarrow \mathbb{N}$  which are (1) bijective (2) maximally compact (3) computationally tractable (i.e., closed form inverses).  $\kappa$  will be used to index  $\{n\}_k^2$ -combinations via the Maculay representation [4] and  $\rho$  will index  $\Sigma^k$  tuples, but is slightly more tricky to define. To maximize compactness, there is an elegant pairing function courtesy of Szudzik [6], which enumerates concentric square shells over the plane  $\mathbb{N}^2$  and can be generalized to hypercubic shells in  $\mathbb{N}^k$ .

Although  $\langle \kappa, \rho \rangle$  could be used directly to exhaustively search the Levenshtein ball, they are temporally biased samplers due to lexicographic ordering.

<sup>2</sup> Following Stirling, we use  $\{n\}_d$  to denote the set of all  $d$ -element subsets of  $\{1, \dots, n\}$ .

Rather, we would prefer a path that uniformly visits every fertile subspace of the Levenshtein ball over time regardless of the grammar or string in question: subsequences of  $\langle \kappa, \rho \rangle$  should discover valid repairs with frequency roughly proportional to the filtration rate, i.e., the density of the admissible set relative to the Levenshtein ball. These additional constraints give rise to two more criteria: (4) ergodicity and (5) periodicity.

To achieve ergodicity, we permute the elements of  $\{n\} \times \Sigma^k$  using a finite field with a characteristic polynomial  $C$  of degree  $m := \lceil \log_p \binom{n}{k} |\Sigma_\varepsilon|^k \rceil$ . By choosing  $C$  to be some irreducible polynomial, one ensures the path has the mixing properties we desire, e.g., suppose  $U : \mathbb{Z}_2^{m \times m}$  is a matrix whose structure is depicted to the right, wherein  $C$  represents a primitive polynomial over  $\mathbb{Z}_2^m$  with known coefficients  $C_{1\dots m}$  and semiring operators  $\oplus := + \pmod{2}$ ,  $\otimes := \wedge$ ,  $\top := 1$ ,  $\circ := 0$ . Since  $C$  is primitive, the sequence  $\mathbf{R} = (U^{0\dots 2^m-1}Y)$  must have *full periodicity*, i.e., for all  $i, j \in [0, 2^m)$ ,  $\mathbf{R}_i = \mathbf{R}_j \Rightarrow i = j$ . To uniformly sample  $\sigma$  without replacement, we construct a partial surjective function from  $\mathbb{Z}_2^m$  onto the Levenshtein ball,  $\mathbb{Z}_2^m \rightarrow \{n\} \times \Sigma_\varepsilon^d$ , cycle over  $\mathbf{R}$ , then discard samples which have no witness in  $\{n\} \times \Sigma_\varepsilon^d$ .

$$U^\top Y = \begin{pmatrix} C_1 & \dots & C_m \\ \top & \circ & \dots & \circ \\ \circ & \dots & \dots & \circ \\ \circ & \dots & \circ & \top & \circ \end{pmatrix}^t \begin{pmatrix} Y_1 \\ \vdots \\ Y_m \end{pmatrix}$$

This procedure requires  $\mathcal{O}(1)$  per sample and roughly  $\binom{n}{d} |\Sigma_\varepsilon|^d$  samples to exhaustively search  $\{n\} \times \Sigma_\varepsilon^d$ . Its acceptance rate  $p^{-m} \binom{n}{d} |\Sigma_\varepsilon|^d$  can be slightly improved with a more suitable base  $p$ , however this introduces some additional complexity and so we elected to defer this optimization.

In addition to its statistically desirable properties, our sampler has the practical benefit of being trivially parallelizable using the leapfrog method, i.e., given  $p$  independent processors, each one can independently check  $[\varphi^{-1}(\langle \kappa, \rho \rangle^{-1}(\mathbf{R}_i), \sigma) \in \mathcal{L}(\mathcal{G})]$  where  $i \equiv p_j \pmod{p}$ . This procedure linearly scales with the number of processors, exhaustively searching  $\Delta_q(\sigma)$  in  $p^{-1}$  of the time required by a single processor, or alternately drawing  $p$  times as many samples in the same time.

To admit variable-length edits and enable deletion, we first define a  $\varepsilon^+$ -production and introduce it to the right- and left-hand side of each terminal unit production contained in the grammar,  $\mathcal{G}$ , as follows:

$$\frac{\mathcal{G} \vdash \varepsilon \in \Sigma}{\mathcal{G} \vdash (\varepsilon^+ \rightarrow \varepsilon \mid \varepsilon \varepsilon^+) \in P} \varepsilon\text{-DUP} \quad \frac{\mathcal{G} \vdash (A \rightarrow B) \in P}{\mathcal{G} \vdash (A \rightarrow B \varepsilon^+ \mid \varepsilon^+ B \mid B) \in P} \varepsilon^+\text{-INT}$$

Finally, to sample  $\sigma \sim \Delta_q(\sigma)$ , we first interleave  $\sigma$  as  $\sigma_\varepsilon$ , then enumerate hole templates  $T(\sigma_\varepsilon, i) = \sigma_{1\dots i-1} \_ \sigma_{i+1\dots n}$  for each  $i \in \cdot \in \{n\}$  and  $d \in 1\dots q$ , then solve for  $A(\sigma_\varepsilon)$ . If  $\sigma := A(\sigma_\varepsilon)$  is nonempty, then each edit from each patch in each  $\tilde{\sigma} \in \sigma$  will match one of the following patterns, covering all three Levenshtein edits:



$$\begin{aligned}
\text{Deletion} &= \left\{ \dots \sigma_{i-1} \text{ } \gamma_1 \text{ } \gamma_2 \text{ } \sigma_{i+1} \dots \mid \gamma_{1,2} = \varepsilon \right\} \\
\text{Substitution} &= \left\{ \begin{aligned} &\dots \sigma_{i-1} \text{ } \gamma_1 \text{ } \gamma_2 \text{ } \sigma_{i+1} \dots \mid \gamma_1 \neq \varepsilon \wedge \gamma_2 = \varepsilon \\ &\dots \sigma_{i-1} \text{ } \gamma_1 \text{ } \gamma_2 \text{ } \sigma_{i+1} \dots \mid \gamma_1 = \varepsilon \wedge \gamma_2 \neq \varepsilon \\ &\dots \sigma_{i-1} \text{ } \gamma_1 \text{ } \gamma_2 \text{ } \sigma_{i+1} \dots \mid \{\gamma_1, \gamma_2\} \cap \{\varepsilon, \sigma_i\} = \emptyset \end{aligned} \right\} \\
\text{Insertion} &= \left\{ \begin{aligned} &\dots \sigma_{i-1} \text{ } \gamma_1 \text{ } \gamma_2 \text{ } \sigma_{i+1} \dots \mid \gamma_1 = \sigma_i \wedge \gamma_2 \notin \{\varepsilon, \sigma_i\} \\ &\dots \sigma_{i-1} \text{ } \gamma_1 \text{ } \gamma_2 \text{ } \sigma_{i+1} \dots \mid \gamma_1 \notin \{\varepsilon, \sigma_i\} \wedge \gamma_2 = \sigma_i \\ &\dots \sigma_{i-1} \text{ } \gamma_1 \text{ } \gamma_2 \text{ } \sigma_{i+1} \dots \mid \gamma_{1,2} = \sigma_i \end{aligned} \right\}
\end{aligned}$$

Although complete with respect to  $\Delta_q(\sigma)$ , this approach can produce patches containing more Levenshtein edits than are strictly necessary to repair  $\sigma$ . To ensure patches are both minimal and syntactically valid, we first introduce a simple technique to minimize the repairs in §1.8. By itself, uniformly sampling minimal repairs  $\tilde{\sigma} \sim \Delta_q(\sigma) \cap \mathcal{L}(\mathcal{G})$  is sufficient but can be quite time-consuming. To further reduce sample complexity and enable real-time repairs, we will then introduce a more efficient density estimator based on adaptive resampling (§1.9).

### 1.8 Patch minimization

Suppose we have a string,  $a ( b$ , and discover the patch,  $\tilde{\sigma} = ( a + b )$ . Although  $\tilde{\sigma}$  is syntactically admissible, it is not minimal. To minimize a patch, we consider the set of all of its constituent subpatches, namely,  $( a + b$ ,  $( a ( b )$ ,  $a + b )$ ,  $( a ( b$ ,  $a + b )$ , and  $a ( b )$ , then retain only the smallest syntactically valid instance(s) by Levenshtein distance. This forms a so-called *patch powerset*, which can be lazily enumerated from the top-down, after which we take all valid strings from the lowest level containing at least one valid string, i.e.,  $a + b$  and  $a ( b )$ . When patches are very large, minimization can be used in tandem with the delta debugging technique [7] to first simplify contiguous edits, then apply the patch powerset construction. Minimization is often useful for estimating the language edit distance: given a single valid repair of arbitrary size, minimization lets us quickly approximate an upper-bound on  $\Delta(\sigma, \ell)$ .

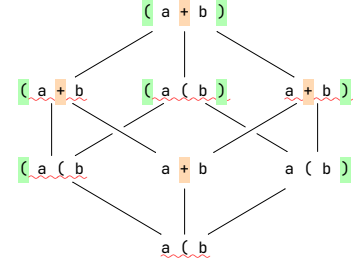


Fig. 1: Patch powerset of  $\tilde{\sigma} = ( a + b )$ .

### 1.9 Probabilistic reachability

Since there are  $\sum_{d=1}^q \binom{n}{d}$  total hole templates, each with  $|\Sigma_\varepsilon|^d$  individual edits to check, if  $n$  and  $q$  are large, this space can be slow to exhaustively search and a

uniform prior may be highly sample-inefficient. Furthermore, naively sampling  $\sigma \sim \Delta_q(\underline{\sigma})$  is likely to produce a large number of unnatural edits and converge poorly on  $\Delta_q(\underline{\sigma}) \cap \mathcal{L}(\mathcal{G})$ . To rapidly rank and render relevant repair recommendations, we prioritize candidate edits according to the following procedure.

(1) Draw samples  $\hat{\sigma} \sim \Delta_q(\underline{\sigma})$  without replacement using §1.7 with leapfrog parallelization. (2) Score by perplexity  $PP(\hat{\sigma})$  using a pretrained variable-order Markov chain (VOMC) [5]. (3) Resample using a concurrent variant of the A-Res [3] online weighted reservoir sampler. (4) Filter Levenshtein edits by admissibility with respect to the grammar, i.e.,  $[\hat{\sigma} \in \mathcal{L}(\mathcal{G})]$ . (5) Minimize and store admissible repairs to a replay buffer,  $\mathcal{Q} \leftarrow \tilde{\sigma}$ , ranked by perplexity. (6) Repeat steps (1)-(5), alternately sampling from the LFSR/VOMC-reweighted online reservoir sampler with probability  $\epsilon$  or stochastically resampled  $\mathcal{Q}$  with probability  $(1 - \epsilon)$ , where  $\epsilon$  decreases from 1 to 0 according to a stepwise schedule relative to the time remaining.

Initially, the replay buffer  $\mathcal{Q}$  is empty and repairs are sampled uniformly without replacement from the Levenshtein ball,  $\Delta_q(\underline{\sigma})$ . As time progresses,  $\mathcal{Q}$  is gradually populated with admissible repairs and resampled with increasing probability, allowing the algorithm to initially explore, then exploit the most promising candidates. This is summarized in Algorithm 1 which is run in parallel across all available CPU cores.

---

**Algorithm 1** Probabilistic reachability

---

**Require:**  $\mathcal{G}$  grammar,  $\underline{\sigma}$  broken string,  $p$  process ID,  $c$  total CPU cores,  $t_{\text{total}}$  timeout.

```

1:  $\mathcal{Q} \leftarrow \emptyset, \mathcal{R} \leftarrow \emptyset, \epsilon \leftarrow 1, i \leftarrow 0, Y \sim \mathbb{Z}_2^m, t_0 \leftarrow t_{\text{now}}$ 
2: repeat
3:   if  $\mathcal{Q} = \emptyset$  or  $\text{Rand}(0, 1) < \epsilon$  then
4:      $\hat{\sigma} \leftarrow \varphi^{-1}(\langle \kappa, \rho \rangle^{-1}(U^{ci+p}Y), \underline{\sigma}), i \leftarrow i + 1$   $\triangleright$  Sample WoR using LFSR.
5:   else
6:      $\hat{\sigma} \sim \mathcal{Q} + \text{Noise}(\mathcal{Q})$   $\triangleright$  Sample replay buffer with additive noise.
7:   end if
8:    $\mathcal{R} \leftarrow \mathcal{R} \cup \{\hat{\sigma}\}$   $\triangleright$  Insert repair candidate  $\hat{\sigma}$  into reservoir  $\mathcal{R}$ .
9:   if  $\mathcal{R}$  is full then
10:     $\hat{\sigma} \leftarrow \text{argmin}_{\hat{\sigma} \in \mathcal{R}} PP(\hat{\sigma})$   $\triangleright$  Select lowest perplexity repair candidate.
11:    if  $\hat{\sigma} \in \mathcal{L}(\mathcal{G})$  then
12:       $\mathcal{Q} \leftarrow \mathcal{Q} \cup \{\hat{\sigma}\}$   $\triangleright$  Insert successful repair into replay buffer.
13:    end if
14:     $\mathcal{R} \leftarrow \mathcal{R} \setminus \{\hat{\sigma}\}$   $\triangleright$  Remove checked sample from the reservoir.
15:  end if
16:   $\epsilon \leftarrow \text{Schedule}((t_{\text{now}} - t_0)/t_{\text{total}})$   $\triangleright$  Update exploration/exploitation rate.
17: until  $t_{\text{total}}$  elapses.
18: return  $\tilde{\sigma} \in \mathcal{Q}$  ranked by  $PP(\tilde{\sigma})$ .
```

---

We would prefer hole templates likely to yield repairs that are (1) admissible (i.e., grammatically correct) and (2) plausible (i.e., likely to have been written by a human author). To do so, we draw holes and rank admissible repairs using a probabilistic distance metric over  $\Delta_q(\sigma)$ .

For example, suppose we are given an invalid string,  $\sigma_\varepsilon : \Sigma^{90}$  and  $\mathcal{Q} \subseteq [0, |\sigma_\varepsilon|) \times \Sigma_\varepsilon^q$ , a distribution over previously successful edits, which we can use to localize admissible repairs. Marginalizing onto  $\sigma_\varepsilon$ , the distribution  $\mathcal{Q}(\sigma_\varepsilon)$  may take the form shown in Fig. 2.

More specifically in our setting, we want to sample from a discrete product space that factorizes into (1) the specific edit locations (e.g., informed by caret position, historical edit locations, or a static analyzer), (2) probable completions (e.g., from a Markov chain or neural language model) and (3) an accompanying *cost model*,  $C : (\Sigma^* \times \Sigma^*) \rightarrow \mathbb{R}$ , which may be any number of suitable distance metrics, such as language edit distance, finger travel distance on a physical keyboard, weighted Levenshtein distance, or stochastic contextual edit distance [2] in the case of probabilistic edits. Our goal then, is to discover repairs which minimize  $C(\sigma, \tilde{\sigma})$ , subject to the given grammar and latency constraints.



Fig. 2: The distribution  $\mathcal{Q}$ , projected onto  $\sigma$ , suggests edit locations likely to yield admissible repairs, from which we draw subsets of size  $d$ .

## References

1. Bringmann, K., Grandoni, F., Saha, B., Williams, V.V.: Truly subcubic algorithms for language edit distance and rna folding via fast bounded-difference min-plus product. *SIAM Journal on Computing* **48**(2), 481–512 (2019)
2. Cotterell, R., Peng, N., Eisner, J.: Stochastic contextual edit distance and probabilistic FSTs. In: *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*. vol. 2 (Short Papers), pp. 625–630. Association for Computational Linguistics, Baltimore, Maryland (June 2014)
3. Efrimidis, P.S.: Weighted random sampling over data streams. *Algorithms, Probability, Networks, and Games: Scientific Papers and Essays Dedicated to Paul G. Spirakis on the Occasion of His 60th Birthday* pp. 183–195 (2015)
4. Knuth, D.E.: *Generating All Combinations and Partitions*, pp. 5–6. Addison-Wesley (2005)
5. Schulz, M.H., Weese, D., Rausch, T., Döring, A., Reinert, K., Vingron, M.: Fast and adaptive variable order markov chain construction. In: *Algorithms in Bioinformatics: 8th International Workshop, WABI 2008, Karlsruhe, Germany, September 15-19, 2008. Proceedings 8*. pp. 306–317. Springer (2008)
6. Szudzik, M.: An elegant pairing function. In: *Special NKS 2006 Wolfram Science Conference*. pp. 1–12 (2006)
7. Zeller, A.: Isolating cause-effect chains from computer programs. *ACM SIGSOFT Software Engineering Notes* **27**(6), 1–10 (2002)