

# Syntax Error Correction as Idempotent Matrix Completion

ANONYMOUS AUTHOR(S)

In this work, we illustrate how to lower context-free language recognition onto a tensor algebra over finite fields. In addition to its theoretical value, this connection has yielded surprisingly useful applications in incremental parsing, code completion and program repair. To accelerate code completion, we design and implement a novel incremental parser-synthesizer that transforms CFGs onto a dynamical system over finite field arithmetic, enabling us to suggest syntax repairs in-between keystrokes. Our approach attains state-of-the-art repair precision@{5, 10} on a dataset of human Python repairs at a fraction of the time and cost of existing methods. Unlike prior work, our approach does not employ neural networks, and is therefore highly flexible to additional constraints and readily extensible to new programming languages.

## 1 INTRODUCTION

In programming, most errors initially manifest as syntax errors, and though often cosmetic, manual repair can present a significant challenge for novice programmers. *Syntax correction* is the problem of repairing a syntactically incorrect program so that it parses. This problem is challenging because well-formed programs have many extra-syntactic constraints, increasing the complexity of repair, but also because the problem is highly under-determined: repairs must be consistent with the surrounding context, and the correct repair is seldom unique, even assuming minimality.

The majority of prior work on syntax correction can be separated into one of two high-level categories: (1) *formal methods*, which use handcrafted rules to identify and correct common syntax errors (e.g., [2]), and (2) *machine learning*, which typically use neural language models to repair code (e.g. [26]). The formal approach can be effective but only produces one repair, while the latter generates more natural edits, but is more costly to train and difficult to incorporate new constraints.

Our primary contribution borrows insights from Boolean matrix parsing [30] and LCL [35] reachability to develop a new theory of error correction based on SAT solving. This approach offers a compelling alternative to its textbook presentation, not only because it unifies parsing, code completion and error correction under a simple algebraic framework, but is highly efficient, composable using ordinary logical primitives and well-suited for integrated developer tools.

In our work, we recast the problem of syntax correction as a special case of tensor completion with semiring algebra, and lay the foundation for a new approach to program repair grounded in formal language theory, unifying *code completion*, *error correction* and *syntax repair*. We provide exact and approximate algorithms for solving these problems in the real-world setting, and implement them in a real-time editor called Tidyparse. Our contribution is a state-of-the-art fast repair procedure that is provably sound and complete up to a Levenshtein bound. Given an arbitrary well-formed grammar, our tool can be used to complete unfinished code, parse incomplete code and repair broken fragments in arbitrary context-free and linear conjunctive languages.

Our paper is structured as follows: we present two high-level approaches to syntax correction: one that samples random edits and accepts only those which parse (Theory 1, probabilistic correction), and another that uses equational reasoning to find the satisfying assignments to a system of multilinear equations over finite fields (Theory 2, model-theoretic). Finally, we show how these two approaches can be combined to attain human-level precision and state-of-the-art wall clock latency on a dataset of human syntax errors and repairs – without using a neural language model.

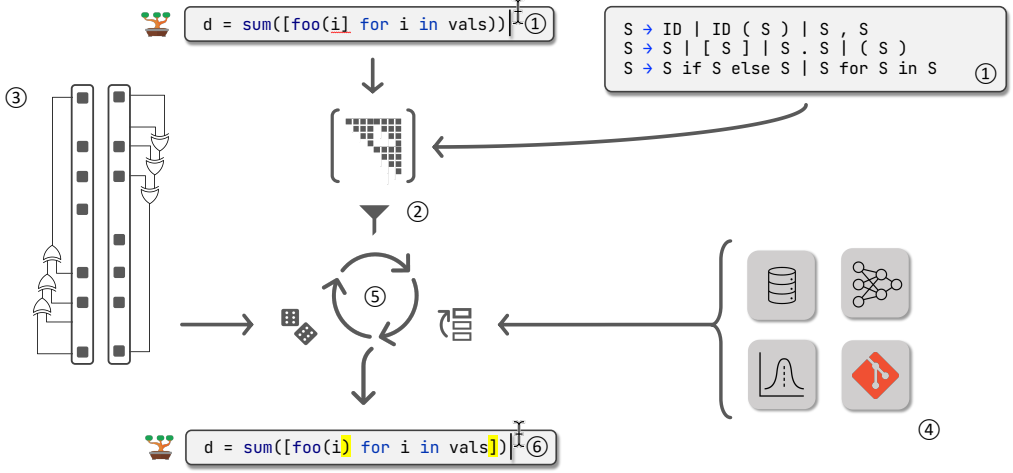


Fig. 1. Our framework consists of three components, namely (2) a solver, (4) ranker and (3) sampler. Given an invalid string and grammar (1), we first compile them into a Boolean tensor (2) representing a multilinear system of equations whose fixed points characterize the admissible set. The system can be solved directly, yielding a set of repairs that are ranked using a suitable scoring function (4). Optionally, we may introduce stochastic edits to the string using the Levenshtein ball sampler (3) and extract the solutions by  $\epsilon$ -decreasing resampling (5). When the language edit distance and solutions are time-sensitive, (5) is typically more efficient.

## 2 OVERVIEW

Our primary contributions in this work are twofold: (1) is a novel reduction from conjunctive language membership with bounded-length string variables to Boolean matrix completion, which we use to develop an algebraic theory of syntax correction, (2) is a randomized algorithm for quickly suggesting syntax repairs to real-world programming errors. We show how to solve this problem exactly and approximately, and implement it in a real-time editor called Tidyparse.

If the string is syntactically valid, we return the parse forest (described in § 5), otherwise, we return the partial parse forest and a small set of edits that would make the string syntactically valid. As we will show in § 3, this set is guaranteed to be nonempty by a straightforward existence proof. To generate repairs, we sample bitvectors without replacement using an LFSR and decode them as repairs. This sequence can be viewed as a randomized space-filling curve over the set of all edits within a fixed distance from the original string. Further details are provided in § 4.7.

Using the Levenshtein ball sampler as a proposal distribution, we then solve for the least fixpoint of an idempotent matrix relation to obtain the set of syntactically admissible repairs, all of which are guaranteed to be minimal and syntactically valid. Finally, these repairs are reranked according to their statistical likelihood, and the most likely samples are presented to the user. We describe our solver in more detail in § 4 and this whole process is illustrated in Fig. 1.

### 2.1 Illustrative example

Below is an example illustrating how our repair procedure works on a real Python code fragment.

- (1) `d = sum([foo(i) for i in vals])`
- (2) 

d	=	sum	(	[	foo	(	i	]	for	i	in	vals	)	)
---	---	-----	---	---	-----	---	---	---	-----	---	----	------	---	---
- (3) 

ID	=	ID	(	[	ID	(	ID	]	for	ID	in	ID	)	)
----	---	----	---	---	----	---	----	---	-----	----	----	----	---	---
- (4) 

ID	=	ID	(	[	ID	(	ID	]	for	ID	in	ID	)	)
----	---	----	---	---	----	---	----	---	-----	----	----	----	---	---

```

(5) 

|    |  |   |    |  |   |   |    |   |    |   |     |    |    |    |   |   |
|----|--|---|----|--|---|---|----|---|----|---|-----|----|----|----|---|---|
|    |  | = |    |  | ( | [ | ID | ( | ID | ] | for | ID | in | ID | ) |   |
| ID |  |   | ID |  | ( | [ |    | ( | ID | ] | for | ID |    | ID | ) | ) |


...
(6) 

|    |   |   |    |   |   |    |   |    |   |     |    |    |    |   |   |
|----|---|---|----|---|---|----|---|----|---|-----|----|----|----|---|---|
| ID | = |   | ID | ( | [ | ID | ( | ID |   | for | ID | in | ID |   | ) |
|    |   | + |    |   |   |    |   |    | ) |     |    |    |    | ] |   |


(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.

Optionally, we may bias the sampler toward previously discovered repairs by adding successful patches to a replay buffer ranked by likelihood. This buffer is then stochastically resampled using a Dirichlet process with stepwise decaying exploration to encourage diversity. As we will show in § 8, this approach is significantly more sample-efficient with respect to language membership and yields much higher precision and repair throughput under low-latency constraints.

## 2.2 Design criteria

During the conception of this work, a number of design choices were made to increase its utility as a real-time programming assistant. To support these design goals, we aim to provide accurate and continuous feedback while the user is typing across a variety of practical programming languages. Our use case motivates a number of optimality criteria, which were taken into account when designing various components and evaluating their performance.

- First and foremost, the framework must be **sound** and **complete**. That is, the parser-generator must (1) accept arbitrary conjunctive grammars, and (2) generate a parser which accepts all and only syntactically valid strings and (3) when a string is invalid, our synthesizer must eventually generate every syntactically valid string within a fixed distance from it.
- Second, we require that the resulting repairs be **plausible** and **diverse**. In other words, the framework should generate repairs that are likely to be written by a human being, consistent with the surrounding context, and reasonably diverse in their structure.
- Third, the framework must be **efficient** and **responsive**. That is, it must be able to recognize well-formed strings in subcubic time, and generate admissible repairs in subpolynomial time. These conditions are necessary to provide real-time feedback whilst the user is typing.
- Fourth, the framework must also be **robust** and **scalable**. In practice, this means that the framework should be robust to multiple errors, handle grammars with a large number of productions and be able to scale linearly as the number of processor cores are increased.
- Finally, the framework must be **flexible** and **extensible**. Intended as a general-purpose tool for generating syntax repairs in a wide variety of programming languages, end-users should be able to extend the framework with their own custom grammars and side-constraints.

Despite the fact that these criteria are not all mutually compatible, we will argue that our framework is able to satisfy most of them simultaneously. In particular, discuss some of the trade-offs made in the design of our framework in § 9.

## 2.3 Research questions

Among the many questions that arise from this work, we focus on the following four:

- (1) **What are the theoretical guarantees?** (e.g. soundness, completeness, optimality)
- (2) **How well does it perform in practice?** (e.g., precision, throughput, wall-clock latency)
- (3) **How does it compare with existing approaches?** (e.g., neural language models)
- (4) **What are the limitations of this technique?** (e.g., how many errors can it handle?)

These questions are answered in the following sections of this paper. We begin by describing the theoretical foundations of our approach in § 4, followed by a description of our parser in § 5. We then evaluate the end-to-end performance of our framework in § 8, and discuss the results in § 9. Finally, we conclude with a summary of our contributions and future work in § 10.

## 3 PROBLEM STATEMENT

The problem of syntax error correction under a finite number of typographic errors is reducible to the bounded Levenshtein CFL reachability problem, which can be formally stated as follows:

*Definition 3.1.* The language edit distance (LED) is 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.,  $\Delta^*(\underline{\sigma}, \ell) := \min_{\sigma \in \ell} \Delta(\underline{\sigma}, \sigma)$ , and  $\Delta$  is the Levenshtein distance.

We seek to find the set of strings  $S$  such that  $\forall \tilde{\sigma} \in S, \Delta(\underline{\sigma}, \tilde{\sigma}) \leq q$ , 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  $\underline{\sigma}$  and denote it  $\Delta_q(\underline{\sigma})$ . Since  $1 \leq \Delta^*(\underline{\sigma}, \ell) \leq q$ , we have  $1 \leq q$ . We now consider an upper bound on  $\Delta^*(\underline{\sigma}, \ell)$ , i.e., the greatest lower bound on  $q$  such that  $\Delta_q(\underline{\sigma}) \cap \ell \neq \emptyset$ .

**LEMMA 3.2.** *For any nonempty language  $\ell : \mathcal{L}$  and invalid string  $\underline{\sigma} : \Sigma^n \cap \bar{\ell}$ , there exists an  $(\tilde{\sigma}, m)$  such that  $\tilde{\sigma} \in \ell \cap \Sigma^m$  and  $0 < \Delta(\underline{\sigma}, \ell) \leq \max(m, n) < \infty$ .*

**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  $\underline{\sigma} \notin \ell$ , it follows that  $0 < \Delta(\underline{\sigma}, \ell)$ . Let us consider two cases, either  $\tilde{\sigma} = \varepsilon$ , or  $0 < |\tilde{\sigma}|$ :

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

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

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

**LEMMA 3.3.** *The complexity of searching  $A$  is upper bounded by  $\mathcal{O}\left(\sum_{c=1}^q \binom{cn+n+1}{c} (|\Sigma| + 1)^c\right)$ .*

**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 size  $n$ . This can be overapproximated by interleaving  $\varepsilon^c$  around every token, i.e.,  $\underline{\sigma}' := (\varepsilon^c \underline{\sigma}_i)_{i=1}^n \varepsilon^c$ , where  $|\underline{\sigma}'| = cn + n + 1$ , 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  $\underline{\sigma}'$  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  $|A| < \sum_{c=1}^q \binom{cn+n+1}{c} |\Sigma_\varepsilon|^c$ .  $\square$

We note that the above bound applies to all strings and languages, and relates to the Hamming bound on  $H_q(\sigma)$ , which only considers substitutions. In practice, much tighter bounds may be obtained by considering the structure of  $\ell$  and  $\sigma$ . For example, based on an empirical evaluation from a dataset of human errors and repairs in Python code snippets ( $|\Sigma| = 50$ ,  $|\sigma| < 40$ ,  $\Delta(\sigma, \ell) \in [1, 3]$ ), we estimate the *filtration ratio*, i.e., the density of the admissible set relative to the Levenshtein ball,  $D = |A|/|\Delta_q(\sigma)|$  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  $q = \Delta^*(\sigma, \ell) \approx 1$ . For example, consider the language of ursine dietary preferences. Although  $\sigma :=$  “Bears like to eat tomatoes” is not in the language, e.g.,  $\tilde{\sigma} :=$  “Bears like to eat” is valid ( $\Delta^* = 1$ ), however there are many other variations with the same or similar edit distance, e.g., or “Bears like to eat {berries, honey, fish}”, or “{Polar, Panda} bears like to eat {seals, bamboo}”. In general, there are many valid strings within a fixed Levenshtein distance of  $\sigma$ , and we seek to find among them, the set of strings that are both syntactically valid and semantically plausible within a short amount of time.

#### 4 MATRIX THEORY

Recall that a CFG is a quadruple consisting of terminals ( $\Sigma$ ), nonterminals ( $V$ ), productions ( $P: V \rightarrow (V \mid \Sigma)^*$ ), and a start symbol, ( $S$ ). It is a well-known fact that every CFG is reducible to *Chomsky Normal Form*,  $P': V \rightarrow (V^2 \mid \Sigma)$ , in which every production takes one of two forms, either  $w \rightarrow xz$ , or  $w \rightarrow t$ , where  $w, x, z : V$  and  $t : \Sigma$ . For example:

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

Given a CFG,  $\mathcal{G}' : \mathbb{G} = \langle \Sigma, V, P, S \rangle$  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,  $\emptyset$  be  $\emptyset$ ,  $\oplus$  be  $\cup$ , and  $\otimes$  be defined as:

$$X \otimes Z := \{ w \mid \langle x, z \rangle \in X \times Z, (w \rightarrow xz) \in P \} \quad (1)$$

If we define  $\sigma_r^\dagger := \{ w \mid (w \rightarrow \sigma_r) \in P \}$ , then construct a matrix with nonterminals on the superdiagonal representing each token,  $M_{r+1=c}^0(\mathcal{G}', e) := \sigma_r^\dagger$  and solve for the fixpoint  $M^* = M + M^2$ ,

$$M^0 := \begin{pmatrix} \emptyset & \sigma_1^\dagger & \emptyset & \dots & \emptyset \\ & \ddots & \ddots & \ddots & \ddots \\ & & \emptyset & \ddots & \emptyset \\ & & & \sigma_n^\dagger & \emptyset \\ \emptyset & \dots & \dots & \dots & \emptyset \end{pmatrix} \Rightarrow \begin{pmatrix} \emptyset & \sigma_1^\dagger & \Lambda & \dots & \emptyset \\ & \ddots & \ddots & \ddots & \ddots \\ & & \emptyset & \ddots & \Lambda \\ & & & \sigma_n^\dagger & \Lambda \\ \emptyset & \dots & \dots & \dots & \emptyset \end{pmatrix} \Rightarrow \dots \Rightarrow M^* = \begin{pmatrix} \emptyset & \sigma_1^\dagger & \Lambda & \dots & \Lambda_\sigma^* \\ & \ddots & \ddots & \ddots & \ddots \\ & & \emptyset & \ddots & \Lambda \\ & & & \sigma_n^\dagger & \Lambda \\ \emptyset & \dots & \dots & \dots & \emptyset \end{pmatrix}$$

we obtain the recognizer,  $R(\mathcal{G}', \sigma) := [S \in \Lambda_\sigma^*] \Leftrightarrow [\sigma \in \mathcal{L}(\mathcal{G})]^1$ .

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,  $\mathbb{1}$ . Note that any encoding which respects linearity  $\varphi(\Lambda \otimes \Lambda') \equiv \varphi(\Lambda) \otimes \varphi(\Lambda')$  is suitable – this particular representation shares the same algebraic structure, but is more widely studied in error correction, and readily compiled into circuits and BLAS primitives. Furthermore, it enjoys the benefit of complexity-theoretic speedups to matrix multiplication.

Details of the bisimilarity between parsing and matrix multiplication can be found in Valiant [30], who first realized its time complexity was subcubic  $\mathcal{O}(n^\omega)$  where  $\omega$  is the asymptotic lower bound for Boolean matrix multiplication ( $\omega < 2.77$ ), and Lee [20], who shows that speedups to CFL parsing were realizable by Boolean matrix multiplication algorithms. While more efficient specialized parsers are known to exist for restricted CFGs, this technique is typically lineararithmic under sparsity and believed to be the most efficient general procedure for CFL parsing.

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

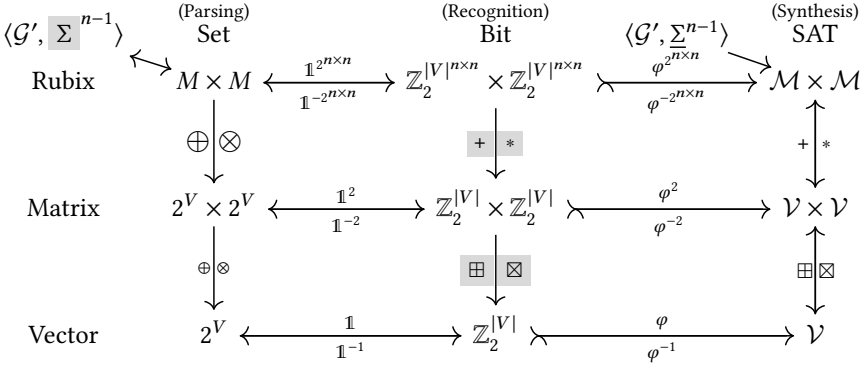
#### 4.1 Parsing with holes

Valiant’s algorithm can be abstracted by lifting it into the domain of linear equations over finite fields, where each bitvector in the upper triangular matrix instead contains polynomials whose solutions identify nonterminals participating in the parse forest of a *porous string* on the superdiagonal (our work). This interpretation of Valiant’s algorithm as an equational theory offers a novel polynomial reduction from language equations onto XORSAT. More formally, we solve the following problem:

**Definition 4.1 (Hole subsumption).** Let  $\underline{\Sigma} := \Sigma \cup \{\_ \}$ , where  $\_$  represents a hole. We denote  $\sqsubseteq: \underline{\Sigma}^n \times \Sigma^n$  as the relation  $\{\langle \sigma, \sigma' \rangle \mid \sigma_i \in \Sigma \implies \sigma_i = \sigma'_i\}$  and the set of all  $\{\sigma' \mid \sigma' \sqsubseteq \sigma\}$  as  $H(\sigma)$ .

**Definition 4.2 (Parsing with holes).** Let  $\underline{\Sigma} := \Sigma \cup \{\_ \}$ , where holes denote free variables over  $\Sigma$ . Given a CFG,  $\mathcal{G} : \mathbb{G}$ , and a *porous string*,  $\sigma : \underline{\Sigma}^{n-1}$ ,  $n : \mathbb{N}^{\geq 3}$ , let us define  $A(\sigma, \mathcal{G}) := H(\sigma) \cap \mathcal{L}(\mathcal{G})$ . We can extract  $A(\sigma, \mathcal{G})$  by encoding  $\langle \sigma, \mathcal{G} \rangle$  as an idempotent matrix,  $\mathcal{M}_\sigma := (\mathbb{Z}_2^{|V|} \rightarrow \mathbb{Z}_2^{|V|})^{n \times n}$ , and solving for its fixedpoints  $[\mathcal{M}_\sigma = \mathcal{M}_\sigma^2]$ , which are then decoded into  $\sigma' \in A(\sigma, \mathcal{G})$ .

The precise encoding is immaterial, so long as  $\langle \boxplus, \boxtimes \rangle$  are defined so the diagram below commutes,<sup>2</sup>



where  $\mathcal{V}$  is a function  $\mathbb{Z}_2^{|V|} \rightarrow \mathbb{Z}_2$ . Note that while always possible to encode  $\mathbb{Z}_2^{|V|} \rightarrow \mathcal{V}$  using the identity function, an arbitrary  $\mathcal{V}$  might have zero, one, or in general, multiple solutions in  $\mathbb{Z}_2^{|V|}$ . In practice, this means that a language equation can be unsatisfiable or underconstrained, however if a solution exists, it can always be decoded into a valid sentence in the language  $\mathcal{L}(\mathcal{G})$ .

So far, we have only considered the syntactic theory of breadth-bounded CFLs with holes, however, our construction can be easily extended to handle the family of CFLs closed under conjunction. The additional expressivity afforded by the language conjunction will be indispensable when considering more practical program repair scenarios that may not be context-free.

#### 4.2 Linear conjunctive reachability

While generally quite expressive, CFLs are themselves not closed under intersection and have other practical limitations, i.e., are unable to express indentation or variable binding. These limitations motivate us towards more expressive yet still efficiently parsable formalisms. In the case of intersection, let us consider the traditional example,  $\mathcal{L}_\cap := \mathcal{L}(\mathcal{G}_1) \cap \mathcal{L}(\mathcal{G}_2)$  defined as follows:

$$\begin{aligned} P_1 &:= \{ S \rightarrow LR, \quad L \rightarrow ab \mid aLb, \quad R \rightarrow c \mid cR \} \\ P_2 &:= \{ S \rightarrow LR, \quad R \rightarrow bc \mid bRc, \quad L \rightarrow a \mid aL \} \end{aligned}$$

<sup>2</sup>Hereinafter, we use gray highlighting to denote types and functions defined over strings and binary constants only.



Note that  $\mathcal{L}_\cap$  generates the language  $\{a^d b^d c^d \mid d > 0\}$ , which according to the pumping lemma is not context-free. However, we can encode the intersection between two or more languages as a single SAT formula by representing each upper-triangular matrix  $\bigcap_{i=1}^c \mathcal{L}(\mathcal{G}_i)$  as a polygonal prism with upper-triangular matrices adjoined to each rectangular face. More precisely, we intersect all terminals  $\Sigma_\cap := \bigcap_{i=1}^c \Sigma_i$ , then for each  $t_\cap \in \Sigma_\cap$  and CFG, construct an equivalence class  $E(t_\cap, \mathcal{G}_i) = \{w_i \mid (w_i \rightarrow t_\cap) \in P_i\}$  and bind them together using conjunction:

$$\bigwedge_{t \in \Sigma_\cap} \bigwedge_{j=1}^{c-1} \bigwedge_{i=1}^{|\sigma|} E(t_\cap, \mathcal{G}_j) \equiv_{\sigma_i} E(t_\cap, \mathcal{G}_{j+1}) \quad (2)$$



Fig. 2. Orientations of a  $\bigcap_{i=1}^4 \mathcal{L}(\mathcal{G}_i) \cap \Sigma^6$  configuration. As  $c \rightarrow \infty$ , this shape approximates a circular cone whose symmetric axis joins  $\sigma_i$  with orthonormal unit productions  $w_i \rightarrow t_\cap$ , and  $[S_i \in \Lambda_\sigma^*]$  represented by the outermost bitvector inhabitants. Equations of this form are equiexpressive with the family of CSLs realizable by finite CFL intersection.

Following Okhotin [22], we then extend our grammar DSL with one additional operator for combining CFGs,  $\wedge : \mathbb{G}^+ \times \mathbb{G}^+ \rightarrow \mathbb{G}^+$ , where  $\mathbb{G}^+$  is a conjunctive grammar (CG). In our setting, CGs naturally subsume CFGs, and in so doing, observe the following denotational semantics:

$$\frac{\Gamma \vdash \mathcal{G} : \mathbb{G}}{\Gamma \vdash \mathcal{G} : \mathbb{G}^+} + \frac{\Gamma \vdash \mathcal{G}, \mathcal{G}' : \mathbb{G}^+}{\Gamma \vdash \mathcal{G} \wedge \mathcal{G}' : \mathbb{G}^+} \wedge \frac{\Gamma \vdash \sigma : \mathcal{L}(\mathcal{G}) \quad \Gamma \vdash \sigma : \mathcal{L}(\mathcal{G}')}{\Gamma \vdash \sigma : \mathcal{L}(\mathcal{G} \wedge \mathcal{G}')} \in \frac{\Gamma \vdash \sigma : \mathcal{L}(\mathcal{G} \wedge \mathcal{G}')}{\Gamma \vdash \sigma : \mathcal{L}(\mathcal{G}) \cap \mathcal{L}(\mathcal{G}')} \cap$$

Given two CFLs  $\mathcal{L}(\mathcal{G}), \mathcal{L}(\mathcal{G}')$ , we can compute their intersection  $\mathcal{L}(\mathcal{G}) \cap \mathcal{L}(\mathcal{G}') \cap \Sigma^d$  by encoding  $[M_{\mathcal{G}, \sigma}^* \equiv_\sigma M_{\mathcal{G}', \sigma}^*]$ . With this feature, we now have the ability to compute the language intersection between two or more CFLs, which will be used to encode bounded Levenshtein CFL reachability.

### 4.3 Levenshtein reachability

The Levenshtein ball is finite and therefore context-free, however materializing this set is intractable for all but the smallest radii and alphabets. Instead, we dynamically instantiate a CFG that can recognize and generate the members of  $\Delta_d(\sigma)$  for any arbitrary  $d : \mathbb{N}, \sigma : \Sigma^*$ . This approach follows from a straightforward extension of Levenshtein automata [27].

In the case where  $\sigma$  and  $\sigma'$  are both fixed strings, Levenshtein distance can be interpreted as a shortest path problem over an unweighted graph whose vertices are the strings  $\sigma, \sigma'$ , and all possible intermediate editor states, and edges represent the Levenshtein edits. Thus viewed, the Levenshtein distance between  $\sigma$  and  $\sigma'$  is simply the geodesic distance [24]. When  $\sigma'$  is instead a free variable and the distance is fixed, we can define a finite automaton accepting all and only strings within Levenshtein distance  $d$  of  $\sigma$  by unrolling the transition dynamics  $\mathcal{L}(\sigma, d)$  up to a fixed horizon  $d$ .

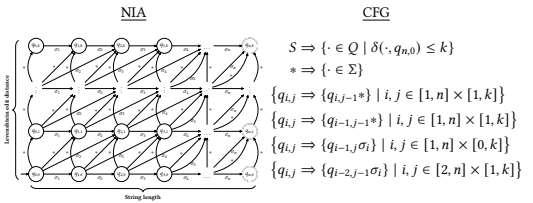


Fig. 3. Levenshtein reachability from  $\Sigma^n$  can be described as an NFA, or a left-linear CFG.

Levenshtein reachability, then, is recognizable by the nondeterministic infinite automaton (NIA) whose topology  $\mathcal{L} = \mathbb{L} \times \mathbb{K}$  factorizes into a product of (a) the monotone Chebyshev topology  $\mathbb{L}$ , equipped with horizontal transitions accepting  $\sigma_i$  and vertical transitions accepting Kleene stars, and (b) the monotone knight's topology  $\mathbb{K}$ , equipped with transitions accepting  $\sigma_{i+2}$ . The structure of this space is approximated by an acyclic NFA, populated by accept states within radius  $k$  of  $q_{n,0}$ , or equivalently, a left-linear CFG whose productions bisimulate the NFA.

#### 4.4 Bounded Levenshtein CFL reachability

Let  $G(\sigma : \Sigma^*, d : \mathbb{N}^+) \mapsto \mathbb{G}$  be the specific construction described in § 4.3 which accepts a string,  $\sigma$ , and an edit distance,  $d$ , and returns a grammar representing the NFA that recognizes the language of all strings within Levenshtein radius  $d$  of  $\sigma$ . To find the language edit distance and corresponding least-distance edits, we must find the smallest  $d$  such that  $\mathcal{L}_d^\cap$  is nonempty, where  $\mathcal{L}_d^\cap$  is defined as  $\mathcal{L}(G(\sigma, d)) \cap \mathcal{L}(\mathcal{G}')$ . In other words, we seek  $\tilde{\sigma}$  and  $d^*$  under which three criteria are all satisfied: (1)  $\tilde{\sigma} \in \mathcal{L}(\mathcal{G}')$ , and (2)  $\Delta(\sigma, \tilde{\sigma}) \leq d^* \iff \tilde{\sigma} \in \mathcal{L}(G(\sigma, d^*))$ , and (3)  $\nexists \sigma' \in \mathcal{L}(\mathcal{G}') . [\Delta(\sigma, \sigma') < d^*]$ . To satisfy these criteria, it suffices to check  $d \in (1, d^*]$  by encoding the Levenshtein automata and the original grammar as a single SAT formula, call it,  $\varphi_d(\cdot)$ , and gradually admitting new acceptance states at increasing radii. If  $\varphi_d(\cdot)$  returns UNSAT,  $d$  is increased until either (1) a satisfying assignment is found or (2)  $d^*$  is attained. Following 3.2, this procedure is guaranteed to terminate in at most either (1) the number of steps required to overwrite every symbol in  $\sigma$ , or (2) the length of the shortest string in  $\mathcal{L}(\mathcal{G}')$ , whichever is greater. When  $\mathcal{L}(\mathcal{G}')$  is context-free,  $\mathcal{L}_d^\cap$  provably context-free using the Bar-Hillel construction [5], however in general, the resulting intersection is a conjunctive language.

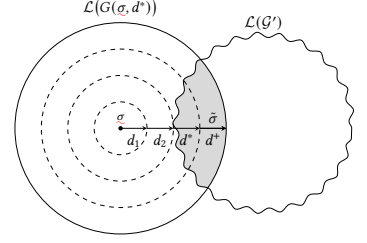


Fig. 4. LED is computed gradually by incrementing  $d$  until  $\mathcal{L}_d^\cap \neq \emptyset$ .

#### 4.5 Tensor sparsification

Although we can encode  $[\mathcal{M}_\sigma = \mathcal{M}_\sigma^2]$  explicitly, naïvely encoding the formula as a dense constraint system is highly suboptimal. For starters, we can always ignore all lower-triangular comparisons  $\{\mathcal{M}_{ij} \mid i \geq j\}$ . A considerably sparser but logically equivalent representation that avoids many elementwise comparisons in the upper-triangular region can be achieved by analyzing the specific string and grammar in question. Since we are typically interested in strings with at least some bound variables,  $\{\sigma : \underline{\Sigma}^* \mid \exists \sigma_i \in \Sigma\}$ , there is likely to be many upper-triangular nonterminals  $\langle v, i, j \rangle \subset V \times \{i, j \mid 0 \leq i \leq j \leq |\sigma|\}$  that are incompatible with  $\sigma$ , i.e., which cannot influence  $[S \in \Lambda_\sigma^*]$  regardless of which values they assume. We call these elements *impossible nonterminals*.

**Definition 4.3 (Impossible nonterminals).** Given  $\sigma : \underline{\Sigma}^*$  and  $\mathcal{G} : \mathbb{G}$ , an impossible nonterminal is a triple  $\langle v, i, j \rangle : V \times \{i, j \mid 0 \leq i \leq j \leq |\sigma|\}$  such that  $\forall \sigma' : \Sigma^{j-i}, \sigma' \in H(\sigma_{i..j}) \implies v \notin \Lambda_{\sigma'}^*$ .

**Definition 4.4 (Terminal closure).** Let  $\Rightarrow_{\mathcal{G}}^*$  be the reflexive transitive closure under  $\mathcal{G}$  defined in the usual way. The terminal closure, denoted  $T(v)$ , is the set  $\{\sigma : \Sigma \mid \exists \alpha, \beta : (V \cup \Sigma)^*, v \Rightarrow_{\mathcal{G}}^* \alpha\sigma\beta\}$ .

**Definition 4.5 (Parikh image).** Let  $\pi : \underline{\Sigma}^* \rightarrow \mathbb{N}^{|\Sigma|}$  be the Parikh vector [23], which counts the number of times each terminal appears in a string. We define the Parikh image as the set of all terminals indexed a nonzero element of the Parikh vector, i.e.,  $\Pi(\sigma : \underline{\Sigma}^*) := \{\sigma' : \Sigma \mid 0 < \pi(\sigma)[\sigma']\}$ .

**LEMMA 4.6.** Given  $\sigma, \mathcal{G}, \langle v, i, j \rangle$ , then  $\Pi(\sigma_{i..j}) \not\subseteq T(v)$  implies  $v \not\Rightarrow^* \sigma_{i..j}$  and  $\langle v, i, j \rangle$  is impossible.

We can overapproximate the set of possible nonterminals by comparing the Parikh image of  $\sigma_{i..j}$  with the terminal closure of  $v$  and eliminating impossible nonterminal triples  $\langle v, i, j \rangle$ . If



$\Pi(\sigma_{i\dots j}) \not\subseteq T(v)$ , then we know  $\sigma_{i\dots j}$  cannot derive  $v$  and  $[(\mathcal{M}_\sigma = \mathcal{M}_\sigma^2)_{i,j,v}]$  may be safely ignored. Note that a tighter approximation can be achieved by considering the minimal Parikh vectors of all  $\sigma : \Sigma^*$  such that  $v \Rightarrow_G^* \sigma \in H(\sigma_{i\dots j})$  rather than just the image on  $V$ , however this is more expensive and  $\Pi(\sigma)$  already discharges a large fraction of impossible nonterminals in practice.

#### 4.6 Isolation

Depicted right is a SAT tensor representing  $\sigma_1 \sigma_2 \sigma_3 \dots$  where shaded regions demarcate known bitvector literals  $\mathcal{L}_{r,c}$  (i.e., representing established nonterminal forests) and unshaded regions correspond to bitvector variables  $\mathcal{V}_{r,c}$  (i.e., representing unknown nonterminal forests to be solved) for an incomplete string. Since  $\mathcal{L}_{r,c}$  are fixed, we precompute them outside the SAT solver.

Clearly, solving complexity is heavily dependent on overall string length. For the sake of complexity, it would be highly advantageous if well-formed subexpressions could be collapsed into a single nonterminal, or multiple nonterminals (in case of ambiguity), to simplify source code snippets without needing to use ad hoc program slicing techniques. Naturally, this raises the question of when can partial derivations be simplified, i.e., under what circumstances is the reduction in Fig. 5 admissible?

This transformation is admissible when the subexpression is “complete”, i.e., its derivation cannot be altered by appending or prepending text. For example, the string  $(-b)$  is morally *complete* in the sense that inserting adjacent text should not alter the interior derivation, whilst  $-b$  is not, as introducing adjacent text (e.g.,  $a-b$ ) may alter the derivation of its contents depending on the structure of the CFG. This question can be reduced to a quotient: Does there exist another nonterminal, when so adjoined that will “strip away” any tokens, leading to another derivation?

More formally, given an arbitrary (potentially ambiguous) context free grammar  $\mathcal{G} : \mathbb{G}$ , and string  $\alpha : \Sigma^*$ , is there a decision procedure that returns whether appending and/or prepending symbols can alter the parse forest of  $\alpha$ ? In other words, we want a function  $F : (\mathbb{G} \times \Sigma^*) \rightarrow \mathbb{B}$  that returns whether  $\alpha$ ’s parse forest according to  $\mathcal{G}$  is unique over  $\beta\alpha\gamma$ , for all  $\beta, \gamma : \Sigma^*$ . Specifically,

**Definition 4.7 (Isolation).** Let  $T_\alpha$  denote the set of all parse trees that are generated by the string  $\alpha$  using  $\mathcal{G}$ , and consider  $T_\alpha$ , the union of all parse trees and their subtrees that (1) can be generated by  $\beta\alpha\gamma$  using  $\mathcal{G}$  for arbitrary  $\beta, \gamma \in \Sigma^*$ , and (2) have a leaf in  $\alpha$ . We call the parse forest  $T_\alpha$  *isolated* iff  $\forall t \in T_\alpha \exists t' \in T_\alpha$ , such that  $t$  is either a subtree of  $t'$ , or  $t'$  is a proper subtree of  $t$ .

If we can identify an  $\alpha : \Sigma$  where  $\alpha = \sigma_{i\dots j}$  for some  $i, j : \mathbb{N}^{<|\sigma|}$ , such that  $\alpha$  has an isolated parse forest, then we can safely replace  $\sigma_{i\dots j} \mapsto \Lambda_\alpha^*$ . For bounded-length strings, this can be solved by padding  $\alpha$  with adjacent holes  $(\_)^n \alpha (\_)^n$  for sufficiently large  $n$ , and checking  $\Lambda_\alpha^\dagger$ , the union of all right- and left-quotients of  $\alpha$  for emptiness, i.e.,  $[\emptyset = \Lambda_\alpha^\dagger \cup \Lambda_\alpha^\dagger]$ . A general solution for arbitrary CFGs,  $\mathcal{G} : \mathbb{G}$  and  $\beta, \gamma \in \Sigma^*$  is more difficult, but would greatly extend the context window and remove the need for ad hoc preprocessing. We leave this as an open problem for future work.

Fig. 5. When is this reduction admissible?

Fig. 6. We can inspect  $\Lambda_\alpha^\dagger, \Lambda_\alpha^\dagger$  to find out.

## 4.7 Sampling without replacement

Now that we have a reliable method to synthesize admissible completions for strings containing holes, i.e., fix *localized* errors,  $S : \mathcal{G} \times (\Sigma \cup \{\varepsilon, \_ \})^n \rightarrow \{\Sigma^n\} \subseteq \mathcal{L}(\mathcal{G})$ , how can we use  $S$  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,  $\sigma \sim \Sigma^{n \pm q} \cap \Delta_q(\sigma)$ , from the Levenshtein  $q$ -ball centered on  $\sigma$ , i.e., the space of all admissible edits with Levenshtein distance  $\leq q$ , loosely analogous to a finite difference approximation over words in a finite language.

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  $\{n\}^2$ -combinations via the Maculay representation [19] 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 [28], which enumerates concentric square shells over the plane  $\mathbb{N}^2$  and can be generalized to hypercubic shells in  $\mathbb{N}^k$ . For our purposes, this generalization will suffice.

Although  $\langle \kappa, \rho \rangle$  could be used directly to exhaustively search the Levenshtein ball, they are temporally biased samplers due to lexicographic ordering. 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: (1) ergodicity and (2) 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  $S = (U^{0 \dots 2^m-1} V)$  must have *full periodicity*, i.e., for all  $i, j \in [0, 2^m)$ ,  $S_i = S_j \Rightarrow i = j$ . To uniformly sample  $\sigma$  without replacement, we construct a partial surjective function from  $\mathbb{Z}_2^m$  onto the Levenshtein ball,  $\mathbb{Z}_2^m \twoheadrightarrow \{n\} \times \Sigma_\varepsilon^d$ , cycle over  $S$ , then discard samples which have no witness in  $\{n\} \times \Sigma_\varepsilon^d$ .

$$U^\top V = \begin{pmatrix} C_1 & \dots & C_m \\ \top & \circ & \circ \\ \circ & \dots & \top \\ \circ & \dots & \circ & \top & \circ \end{pmatrix}^t \begin{pmatrix} V_1 \\ \vdots \\ V_m \end{pmatrix}$$

This procedure requires  $\tilde{O}(1)$  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  $[\varphi^{-1}(\langle \kappa, \rho \rangle^{-1}(S_i), \sigma) \in \mathcal{L}(\mathcal{G})]$  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 amount of time.

To admit variable-length edits, we first define a  $\varepsilon^+$ -production and introduce it to the right- and left-hand side of each terminal in a unit production:

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

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

Finally, to sample  $\sigma \sim \Delta_q(\sigma)$ , we enumerate templates  $H(\sigma, i) = \sigma_{1\dots i-1} \_ \sigma_{i+1\dots n}$  for each  $i \in \cdot \in \{d\}^n$  and  $d \in 1 \dots q$ , then solve for  $\mathcal{M}_\sigma^*$ . If  $[S \in \Lambda_\sigma^*]$ , then each edit in each  $\tilde{\sigma} \in \sigma$  will match one of the following seven patterns:

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

This approach is tractable for  $n \lesssim 100, q \lesssim 3$ , however more complex repairs require a more efficient density estimator. We will now discuss an approach that uses an adaptive sampler (§ 4.8).

#### 4.8 Probabilistic reachability

Since there are  $\sum_{d=1}^q \binom{n}{d}$  total sketch templates, each with  $(|\Sigma| + 1)^d$  individual edits to check, if  $n$  and  $q$  are large, this space can be intractable to exhaustively search and a uniform prior may be highly sample-inefficient. Furthermore, naïvely sampling  $\sigma_i \sim \Sigma^{n \pm q} \cap \Delta_q(\sigma)$  is likely to produce a large number of unnatural edits. To provide rapid and relevant suggestions, we prioritize candidate repairs according to the following seven-step procedure:

- (1) Retrieve the most recent grammar,  $\mathcal{G}$ , and string,  $\sigma$ , from the editor.
- (2) Draw samples  $\sigma_i \sim \{d\}^n \times \Sigma^{n \pm q} \cap \Delta_q(\sigma)$  without replacement using Eq. 4.7.
- (3) Filter completions by admissibility with respect to the grammar,  $\sigma_i \cap \mathcal{L}_\mathcal{G}$ .
- (4) Rerank admissible repairs by the edit cost model,  $C(\sigma, \tilde{\sigma})$ .
- (5) Display the top-k repairs by edit cost found within p-seconds to the user.

Suppose we are given an invalid string,  $\sigma : \Sigma^{90}$  and  $\mathcal{F}_\theta$ , a distribution over possible edits locations provided by a probabilistic or neural language model, which we can use to localize admissible repairs. For example, by marginalizing onto  $\sigma$ , the distribution  $\mathcal{F}_\theta$  could take the following form:

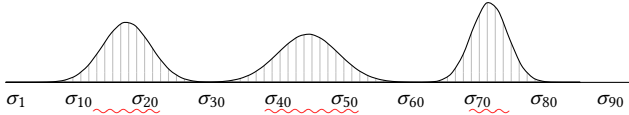


Fig. 7. The distribution  $\int \mathcal{F}_\theta(\cdot \mid \sigma) d\theta$ , projected onto the invalid string, suggests edit locations most likely to yield admissible repairs, from which we draw subsets of size  $d$ .

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 [12] 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.

#### 4.9 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,  $( a + b$ ,  $( a ( b )$ ,  $a + b )$ ,  $( a ( b, a + b )$ ,  $( a ( b, a + b$ , 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 using the Maculay representation, after which we take all valid elements from the lowest level containing at least one admissible element, i.e.,  $a + b$  and  $a ( b )$ . When patches are very large, minimization can be used in tandem with the delta debugging technique [34] 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  $\Delta(\sigma, \ell)$  – much tighter than indicated by Lemma 3.2.

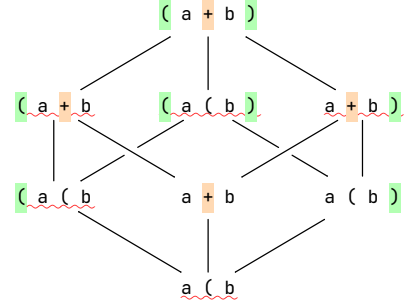


Fig. 8. Patch powerset of  $\tilde{\sigma} = ( a + b )$ .

## 5 PARSING

Although parsing is not the primary objective of this work, it is an integral component of the repair process, and one of the advantages of using a matrix representation is its excellent error recovery properties. Unlike traditional parsers which fail on a error, matrix-based parsers support parsing invalid strings. In this section, we will describe how to parse strings using matrix powering, and how those results inform the repair process.

### 5.1 Tree Denormalization

Our parser emits a binary forest consisting of parse trees for the candidate string which are constructed bottom-up using a variant of  $\otimes$  called  $\hat{\otimes}$ , which simply records backpointers:

$$X \hat{\otimes} Z := \left\{ w \begin{matrix} \nearrow x \\ \searrow z \end{matrix} \mid \langle x, z \rangle \in X \times Z, (w \rightarrow xz) \in P \right\} \quad (3)$$

Due to Chomsky normalization however, the resulting forests are full of trees that are thin and crooked. To restore the natural shape of the tree, we first construct the parse forests bottom-up, then prune away synthetic nonterminals top-down by recursively grafting denormalized grandchildren onto the root. This transformation is purely cosmetic and only used when rendering the parse trees.

**Algorithm 1** Tree denormalization

---

```

procedure CUT( $t$ : Tree)
  stems  $\leftarrow \{\text{CUT}(c) \mid c \in t.\text{children}\}$ 
  if  $t.\text{root} \in (V_{\mathcal{G}'} \setminus V_{\mathcal{G}})$  then
    return stems
  else
    return  $\{\text{Tree}(t.\text{root}, \text{stems})\}$ 
  end if
end procedure

```

---

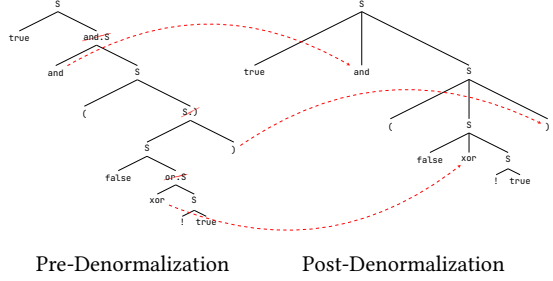


Fig. 9. Since  $\mathcal{G}'$  contains synthetic nodes, to recover a parse tree congruent with the original grammar  $\mathcal{G}$ , we prune all synthetic nodes and graft their stems onto the grandparent via a simple recursive procedure (Alg. 1).

## 5.2 Relation between parsing and repair

Parsing and repair are intimately related: parsing tells us roughly where the repairs should occur and repair is aided by the results of parsing; parseable subtrees can be used to guide the repair process. Although reparsing may be avoided after repair, it is not straightforward how to decode the solution to the matrix equivalence relation as a parse forest, so we resort to decoding the generating string, then reparsing to obtain the CSTs.

The process of parsing the string  $S + S = S$  by matrix powering is depicted below. This occurs by repeatedly joining rooted subtrees using the  $\hat{\otimes}$  operator. In this particular case, the parse forest is unique, although in general,  $M_{ij}$  may contain multiple subtrees should ambiguity arise:

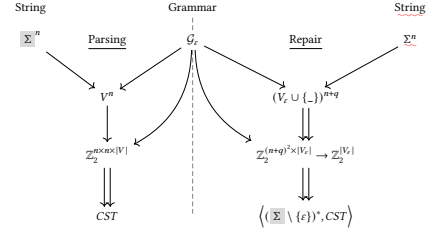


Fig. 10. Parsing and repair share the same grammar representation, i.e., LCGs.



Fig. 11. The parse forest for the string  $S + S = S$  is constructed incrementally by computing  $M_{\sigma}^2 + M_{\sigma} = M_{\sigma}$ .

Our algorithm produces set of concrete syntax trees (CSTs) for a given valid string. Otherwise, if the string is invalid, the algorithm generates a set of admissible corrections, alongside their CSTs. The fixpoint search is sure to converge in at most  $|\sigma| - 1$  steps, as each iteration will solve for a single diagonal of the matrix, each entry of which may contain up to  $|V_{\epsilon}|$  rooted subtrees.

In the case of repair where entries are SAT variables, we can encode the equivalence relation  $\varphi[M_{\sigma} = M_{\sigma}^2]$  directly, instead of solving the fixed point iteratively. Instead of trying to decode the trees directly from the SAT solver, the satisfying models are decoded as strings then fed back into the parser to obtain the parse forest, which is guaranteed to exist if the system is satisfiable.

### 5.3 Incrementalization

When a parsed string is altered, we can reuse prior work by only recomputing affected submatrices, yielding a reparser whose complexity is location-dependent, i.e., at worst quadratic in terms of  $|\Sigma^*|$  assuming  $\mathcal{O}(1)$  cost for each CNF-nonterminal subset join,  $V'_1 \otimes V'_2$ . Letting shaded nodes represent observed or bound variables and unshaded nodes as free variables, we depict the worst-case post-editing state of the parse trellis in Fig. 12.

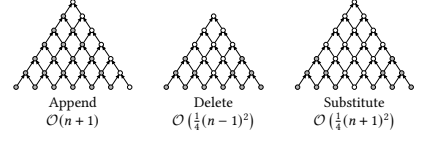


Fig. 12. Incremental reparsing only requires computing submatrices affected by an edit.

The problem of incremental parsing is closely related to *dynamic matrix inversion* in the linear algebra setting, and *incremental transitive closure* with vertex updates in the graph setting. By carefully encoding the matrix relation from § 4 and employing an incremental SAT solver, we can gradually update SAT constraints as new keystrokes are received to eliminate redundancy.

We can use an incremental SAT solver to encode the constraints. To do so, we encode the matrix equivalence relation and use the incremental SAT solver to update the constraints as new keystrokes are received. The incremental SAT solver will eliminate redundancy as it is introduced. Once a solution is found, we can use the algorithm from § 5.1 to recover the parse tree, then introduce a new constraint to block the solution and continue searching for other solutions.

### 5.4 Error Recovery

Not only is Tidyparse capable of suggesting repairs to invalid strings, it can also return partial trees for those same strings, which is often helpful for debugging purposes. Unlike LL- and LR-style parsers which require special rules for error recovery, Tidyparse can simply analyze the structure of  $M^*$  to recover parse branches. If  $S \notin \Lambda_\sigma^*$ , the upper triangular entries of  $M^*$  will take the form of a jagged-shaped ridge whose peaks signify the roots of maximally-parsable substrings  $\hat{\sigma}_{i,j}$ .

These branches are located on peaks of the upper triangular (UT) matrix ridge. As depicted in Fig. 14, we traverse the peaks by decreasing elevation to collect partial AST branches and display the highest nonoverlapping branches, in this case  $T_C$  and  $T_A$  to the user, to help them diagnose the parsing error and manually repair it.

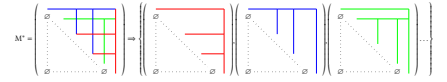


Fig. 13. The matrix  $M^*$  contains all admissible binary trees of a fixed breadth.

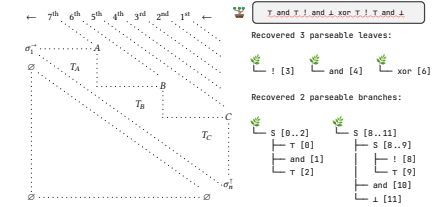


Fig. 14. We can recover the partial subtrees for invalid strings by inspecting the  $M^*$ .

### 5.5 Nonterminal Stubs

Tidyparse augments CFGs with two additional rules, which are desugared into a vanilla CFG before parsing. The first rule,  $\alpha$ -SUB, allows the user to define a nonterminal parameterized by  $\alpha$ , a non-recursive nonterminal in the same the CFG representing some finite type and its inhabitants.  $\alpha$ -SUB replaces all productions containing  $\langle \alpha \rangle$  with the terminals in their transitive closure,  $\alpha \rightarrow^* \beta$ . The second rule,  $\alpha$ -INT, introduces homonymous terminals for each user-defined nonterminal.

$$\frac{\mathcal{G} \vdash (w\langle \alpha \rangle \rightarrow xz) \in P \quad \alpha^* : \{\beta \mid (\alpha \rightarrow^* \beta) \in P\}}{\mathcal{G} \vdash \forall \beta \in \alpha^*. (w\langle \alpha \rangle \rightarrow xz) [\beta/\alpha] \in P'} \quad \alpha\text{-SUB} \quad \frac{\mathcal{G} \vdash v \in V}{\mathcal{G} \vdash (v \rightarrow \langle v \rangle) \in P} \quad \langle \cdot \rangle\text{-INT}$$



Some dependently typed programming languages can do parsing in the type checker – our parser can also perform some form of type checking. Tidyparse automatically expands typed expressions into ordinary nonterminals using the  $\alpha$ -SUB rule, for example when parsing an expression of the form  $x + y$ , the grammar will recognize `true + false` and `1 + 2`, but not `1 + true`.



```
E<X> -> E<X> + E<X> | E<X> * E<X> | ( E<X> )
X -> Int | Bool

E<Int> -> E<Int> + E<Int> | E<Int> * E<Int>
E<Bool> -> E<Bool> + E<Bool> | E<Bool> * E<Bool>
```

When completing a bounded-width string, one finds it is often convenient to admit nonterminal stubs, representing unexpanded subexpressions. To enable this functionality, we introduce a synthetic production for each  $v \in V$  using the  $\langle \cdot \rangle$ -INT rule. Users can interactively build up a complex expression by placing the caret over a stub, then pressing `ctrl` + `Space`:



```
false or ! true or <S> and <S> or <S>

1.) false or ! true or true and <S> or <S>
2.) false or ! true or false and <S> or <S>
3.) false or ! true or ! <S> and <S> or <S>
4.) false or ! true or <S> and <S> and <S> or <S>
5.) false or ! true or <S> or <S> and <S> or <S>
...
```

This functionality can also be useful inside a completion, which might be expanded as follows:



```
if <Vexp> _ _ _ _

1.) if map X then <Vexp> else <Vexp>
2.) if uncurry X then <Vexp> else <Vexp>
3.) if foldright X then <Vexp> else <Vexp>
...
```

## 6 PRACTICAL EXAMPLE

Suppose we are given the following context-free grammar:



```
S -> S and S | S xor S | ( S ) | true | false | ! S
```

For reasons that will become clear in Sec. 4, this will automatically be rewritten into the grammar:

```
F.! -> !      S.) -> S F.) and.S -> F.and S   S -> F.! S   S -> false   S -> S ε+
F.( -> (      F.xor -> xor   xor.S -> F.xor S   S -> S and.S   S -> true     ε+ -> ε
F.) -> )      F.and -> and    S -> S xor.S   S -> F.( S.)   S -> <S>     ε+ -> ε+ ε+
```

Given a string containing holes, our tool will return several completions in a few milliseconds:



```

true _ _ _ ( false _ ( _ _ _ _ ! _ _ ) _ _ _ _
true xor ! ( false xor ( <S> ) or ! <S> ) xor <S>
true xor ! ( false and ( <S> ) or ! <S> ) xor <S>
true xor ! ( false and ( <S> ) and ! <S> ) xor <S>
true xor ! ( false and ( <S> ) and ! <S> ) and <S>
...

```

Similarly, if provided with a string containing various errors, it will return several suggestions how to fix it, where **green** is insertion, **orange** is substitution and **red** is deletion.



```

true and ( false or and true false
1.) true and ( false or ! true )
2.) true and ( false or <S> and true )
3.) true and ( false or ( true ) )
...
9.) true and ( false or ! <S> ) and true false

```

Given an invalid string, the tool will first map the string to a coarsened version and generate edits:



```

v = [float(n for n in l.split(':'))]
v = [float(n for n in l.split(':'))]
v = [float()n for n in l.split(':')]
v = [float() for n in l.split(':')]
v = [float(n) for n in l.split(':')]

```



```

w = [ w ( w w w w w . w ( w ) ]
w = [ w ( w w w w w . w ( w ) ]
w = [ w ( ) w w w w w . w ( w ) ]
w = [ w ( ) w w w w w . w ( w ) ]
w = [ w ( w ) w w w w w . w ( w ) ]

```

This coarsening is done to reduce the number of possible corrections, and is admissible because CFLs are closed under homomorphisms. (If we wanted to provide a lexical expansion, this would be possible.) These candidates then reranked using a probability metric.



```

v = [float(n for n in l.split(':'))]
v = [float(n) for n in l.split(':')]
v = [float(n for n in l.split(':'))]
v = [float()n for n in l.split(':')]
v = [float() for n in l.split(':')]

```



```

w = [ w ( w w w w w . w ( w ) ]
w = [ w ( w ) w w w w w . w ( w ) ]
w = [ w ( w w w w w . w ( w ) ]
w = [ w ( ) w w w w w . w ( w ) ]
w = [ w ( ) w w w w w . w ( w ) ]

```

Finally, we eliminate everything which is not parsed by the official Python parser.



```
v = [float(n for n in l.split(':'))]
-----34-----
v = [float(n) for n in l.split(':')]
-----34-----
v = [float() for n in l.split(':')]
```



```
w = [ w ( w w w w w . w ( w ) ) ]
-----34-----
w = [ w ( w ) w w w w . w ( w ) ]
-----34-----
w = [ w ( ) w w w w . w ( w ) ]
```

Tidyparse requires a grammar, which can either be provided by the user or ingested from a BNF specification. The following grammar represents a slightly more realistic programming language:



```
S -> A | V | ( X , X ) | X X | ( X )
A -> Fun | F | L | L in X
Fun -> fun V -> X
F -> if X then X else X
L -> let V = X | let rec V = X
V -> Vexp | ( Vexp ) | Vexp Vexp
Vexp -> VarName | FunName | Vexp V0 Vexp | ( VarName , VarName ) | Vexp Vexp
VarName -> a | b | c | d | e | ... | z
FunName -> foldright | map | filter | curry | uncurry
V0 -> + | - | * | / | > | = | < | `| | &&
---
let curry f = ( fun x y -> f ( _ _ ) )
-----
let curry f = ( fun x y -> f ( <X> ) )
let curry f = ( fun x y -> f ( <FunName> ) )
let curry f = ( fun x y -> f ( curry <X> ) )
...
```

We can also handle error correction and completion in the untyped  $\lambda$ -calculus, as shown below:



```
sxp -> λ var . sxp | sxp sxp | var | ( sxp ) | const
const -> 1 | 2 | 3 | 4 | 5 | 6
var -> a | b | c | f | x | y | z
---
( λ f . ( λ x . f ( x x ) ) ( λ x . f ( x x ) )
-----
1.) ( λ f . ( λ x . f ( x x ) ) ) λ x . f ( x x )
2.) ( λ f . ( λ x . f ( x x ) ) x ) λ x . f ( x x )
3.) ( λ f . ( λ x . f ( x x ) ) ( λ x . f ( x ) ) )
...
```

Name resolution and scope checking is also possible but requires a more sophisticated grammar.

## 6.1 Grammar Assistance

Tidyparse uses a CFG to parse the CFG, so it can provide assistance while the user is designing the CFG. For example, if the CFG does not parse, it will suggest possible fixes. In the future, we intend to use this functionality to perform example-based codesign and grammar induction.



```
B -> true | false | 
-----
B -> true | false 
B -> true | false <RHS>
B -> true | false | <RHS>
...
```

## 6.2 Interactive Nonterminal Expansion

Users can interactively build up a complex expression by placing the caret over a stub they wish to expand, then pressing `ctrl` + `Space`:



```
if <Vexp> X then <Vexp> else <Vexp>
-----
if map X then <Vexp> else <Vexp>
if uncurry X then <Vexp> else <Vexp>
if foldright X then <Vexp> else <Vexp>
...
```

## 6.3 Conjunctive Grammars

Many natural and programming languages exhibit context-sensitivity, such as Python indentation. Unlike traditional parser-generators, Tidyparse can encode CFL intersection, allowing it to detect and correct errors in a more expressive family of languages than would ordinarily be possible using CFGs alone. For example, consider the grammar from Sec. 4.2:



```
S -> L R    L -> a b | a L b    R -> c | c R    &&&
S -> L R    R -> b c | b R c    L -> a | a L
---
- - - - -
1.) a b c
2.) a a b b c c
3.) a a a b b b c c c
4.) a a a a b b b b c c c c
...
```

Tidyparse uses the notation  $G_1 \&\&\& G_2$  and  $G_1 ||| G_2$  to signify  $\mathcal{L}_{G_1} \cap \mathcal{L}_{G_2}$  and  $\mathcal{L}_{G_1} \cup \mathcal{L}_{G_2}$  respectively, i.e., the intersection or union of two or more grammars' languages. Composition, complementation and other operations on finite languages are also possible, although undocumented at present.

## 7 RELATED WORK

Three important questions arise when repairing syntax errors: (1) is the program broken in the first place? (2) if so, where are the errors located? (3) how should those locations then be altered? In the case of syntax correction, those questions are addressed by three related research areas, (1) parsing, (2) language equations and (3) repair. We survey each of those areas in turn.

### 7.1 CFL Parsing

Context-free language (CFL) parsing is the well-studied problem of how to turn a string into a unique tree, with many different algorithms and implementations (e.g., shift-reduce, recursive-descent, LR). Many of those algorithms expect grammars to be expressed in a certain form (e.g., left- or right- recursive) or are optimized for a narrow class of grammars (e.g., regular, linear).

General CFL parsing allows ambiguity (non-unique trees) and be formulated as a dynamic programming problem, as shown by Cocke-Younger-Kasami (CYK) [25], Earley [15] and others. These parsers have roughly cubic time complexity with respect to the length of the input string.

As shown by Valiant [30], Lee [20] and others, general CFL recognition is in some sense equivalent to binary matrix multiplication, another well-studied combinatorial problem with broad applications, known to be at worst subcubic. This realization unlocks the door to a wide range of complexity-theoretic and practical speedups to CFL recognition and fast general parsing algorithms.

### 7.2 Language Equations

Language equations are a powerful tool for reasoning about formal languages and their inhabitants. First proposed by Ginsburg et al. [16] for the ALGOL language, language equations are essentially systems of inequalities with variables representing *holes*, i.e., unknown values, in the language or grammar. Solutions to these equations can be obtained using various fixpoint techniques, yielding members of the language. This insight reveals the true algebraic nature of CFLs and their cousins.

Being an algebraic formalism, language equations naturally give rise to a kind of calculus, vaguely reminiscent of Leibniz' and Newton's. First studied by Brzozowski [9, 10] and Antimirov [3], one can take the derivative of a language equation, yielding another equation, which can be interpreted as a kind of continuation or language quotient, returning the suffixes that complete a given prefix. This technique leads to an elegant family of algorithms for incremental parsing [1, 21] and automata minimization [8]. In our setting, differentiation corresponds to code completion.

In this paper, we restrict our attention to language equations over context-free and weakly context-sensitive languages, whose variables coincide with edit locations in the source code of a computer program, and solutions correspond to syntax repairs. Although prior work has studied the use of language equations for parsing [21], to our knowledge they have never previously been considered for the purpose of code completion or syntax error correction.

### 7.3 Syntax Repair

In finite languages, syntax repair corresponds to spelling correction, a more restrictive and largely solved problem. Schulz and Stoyan [27] construct a finite automaton that returns the nearest dictionary entry by Levenshtein edit distance. Though considerably simpler than syntax correction, their work shares similar challenges and offers insights for handling more general repair scenarios.

When a sentence is grammatically invalid, parsing grows more challenging. Like spelling, the problem is to find the minimum number of edits required to transform an arbitrary string into a syntactically valid one, where validity is defined as containment in a (typically) context-free language. Early work, including Aho [2, 17] propose a dynamic programming algorithm to computes the minimum number of edits required to fix invalid string. Prior work on error correcting parsing

only considers the shortest edit(s), and does not study multiple edits or the Levenshtein ball, in which the problem of actually generating the repairs is not well-posed, as there are often many valid strings that can be obtained within a given number of edits. We instead focus on bounded Levenshtein reachability, which is the problem of finding useful repairs within a fixed Levenshtein distance of the broken string, which requires language intersection.

#### 7.4 Program Synthesis

There is related work on string constraint solving in the constraint programming literature, featuring solvers like CFGAnalyzer and HAMPI [18], which consider bounded context free grammars and intersections thereof. Axelson et al. (2008) [4] has some work on incremental SAT encoding but does not exploit the linear-algebraic structure of parsing nor provide real-time guarantees. Finally, Loris D’Antoni did some great work on *symbolic automata* [13], a generalization of finite automata which allow infinite alphabets and symbolic expressions over them. In none of the constraint programming literature we surveyed do any of the approaches employ matrix-based parsing, and therefore do not enjoy the optimality guarantees of Valiant’s parser. Our solver can handle CFGs and conjunctive grammars with finite alphabets and does not require any special grammar encoding. The matrix encoding makes it particularly amenable to parallelization.

#### 7.5 Error Correction

Our work focuses on errors arising from human factors in computer programming, in particular *syntax error correction*, which is the problem of fixing partially corrupted programs. Modern research on error correction, however, can be traced back to the early days of coding theory when researchers designed *error-correcting codes* (ECCs) to denoise transmission errors induced by external interference, e.g., collision with a high-energy proton, manipulation by an adversary or even typographical mistake. In this context, *code* can be any logical representation for communicating information between two parties (such as a human and a computer), and an ECC is a carefully-designed which ensures that even if some portion of the message should become corrupted, one can still recover the original message by solving a linear system of equations. When designing ECCs, one typically assumes a noise model over a certain event space, such as the Hamming [14, 29] or Levenshtein [6, 7] balls, from which we draw inspiration for our work.

### 8 EXPERIMENTAL SETUP

To evaluate our model, we primarily use two datasets. We first evaluate it on 5,600 pairs of (broken, fixed) Python code snippets from Wong et al.’s StackOverflow dataset [32] shorter than 40 lexical tokens, whose minimized patch sizes shorter than four lexical tokens ( $|\Sigma| = 50$ ,  $|\sigma| < 40$ ,  $\Delta(\sigma, \ell) < 4$ ). Minimization uses the Delta debugging [34] technique to isolate the smallest lexical patch that repairs a broken Python snippet.

In the first set of experiments, we uniformly sample without replacement from the Levenshtein edit ball using a LFSR over  $\mathbb{Z}_2^n$ , and measure the Precision@k of our repair procedure against human repairs of varying edit distances and latency cutoffs. This provides a baseline for the relative density of the admissible set, and an upper bound on the latency of attaining a given precision.

In the second set of experiments, we use an adaptive sampler that stochastically resamples edits using a Dirichlet process. Repairs are scored and placed into a ranked buffer, from which new edits are resampled with frequency relative to their perplexity, and additive noise is introduced. Initially, the sampler draws uniformly but is gradually biased towards more likely repairs.

To train the adaptive sampler we use a length-5 variable-order Markov (VOM) chain on 55 million StackOverflow tokens from the BIFI [33] dataset, which takes roughly 10 minutes on 40 Intel Skylake cores running at 2.4 GHz. Sequences are scored using negative log likelihood with

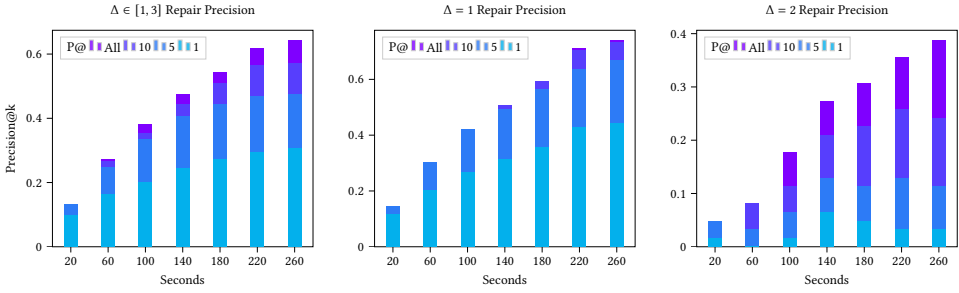


Laplace smoothing, or Levenshtein edits when no dataset is available. Evaluation measures the Precision@{1, 5, 10, All} for samples at varying latency cutoffs.

Both sets of experiments use 40 Intel Skylake cores, running at 2.4 GHz, with 16 GB of RAM. We measure the precision on an exact lexical match abstract repair, following the Seq2Parse [26] evaluation, and give a baseline for their approach on the same dataset. Unlike Seq2Parse, our approach does not require a training corpus, and can be used to repair arbitrary context-free and linear conjunctive languages.

### 8.1 Uniform Sampling benchmark

Below, we plot the results of a human repair benchmark measuring the Precision@k of our repair procedure against human repairs of varying edit distances and latency cutoffs.



For comparison below are the results from Seq2Parse on the same dataset. Seq2Parse only supports Precision@1 repairs, and so we only report precision@1 from the StackOverflow benchmark for comparison. Unlike our approach which only produces syntactically correct repairs, Seq2Parse also produces syntactically incorrect repairs and so we report the percentage of syntactically correct repairs (Syntactic), as well as the precision of the abstract tokens repairs (AbstractEval), and the exact character match precision (CharMatch).

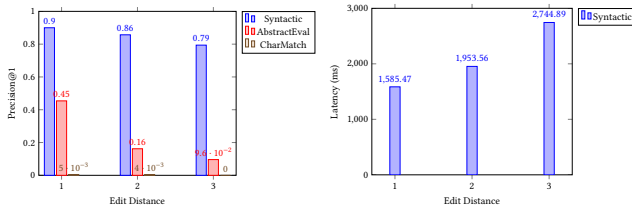


Fig. 15. Seq2Parse precision@1 and latency on the StackOverflow dataset.

### 8.2 Repair with an adaptive sampler

In the following benchmark, we measure the Precision@k of our repair procedure against human repairs of varying edit distances and latency cutoffs, using an adaptive resampling procedure. This sampler maintains a buffer of successful repairs ranked by perplexity and uses stochastic local search to resample edits within a neighborhood. Initially, edits are sampled uniformly at random. Over time and as the admissible set grows, it prioritizes edits nearby low-perplexity repairs. This technique offers a significant advantage in the low-latency setting.

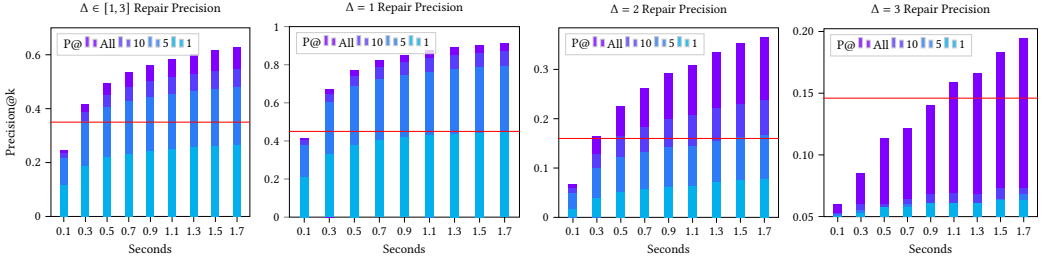
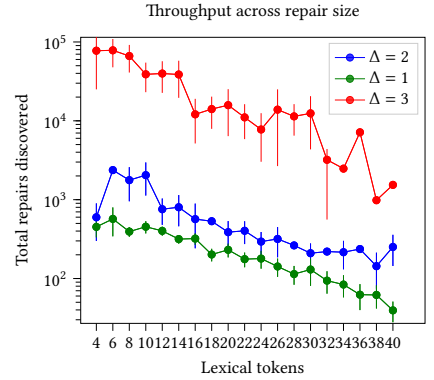


Fig. 16. Adaptive sampling repairs. The red line indicates Seq2Parse precision@1 on the same dataset. Since it only supports generating one repair, we do not report precision@k or the intermediate latency cutoffs.

### 8.3 Throughput benchmark

We evaluate throughput by sampling edits across invalid strings of varying length, measuring the total number of syntactically valid strings discovered, and the human patch size. Each trial is terminated after 10 seconds, and the experiment is repeated across 7.3k total repairs.

While the sampler is able to saturate the admissible set for 1- and 2-edit repairs before the timeout elapses, 3-edit throughput is heavily constrained by compute around 16 lexical tokens, when Python’s Levenshtein ball has a volume of roughly  $6 \times 10^8$  edits. This bottleneck can be relaxed with a longer timeout or additional CPU cores. Despite the high computational cost of sampling multi-edit repairs, our Precision@All remains competitive with the Seq2Parse baseline at the same latency.



### 8.4 Synthetic repair benchmark

In the following benchmarks, we measure the wall clock time required to synthesize solutions to length-50 strings sampled from various Dyck languages, where Dyck-n is the Dyck language containing n different types of balanced parentheses.

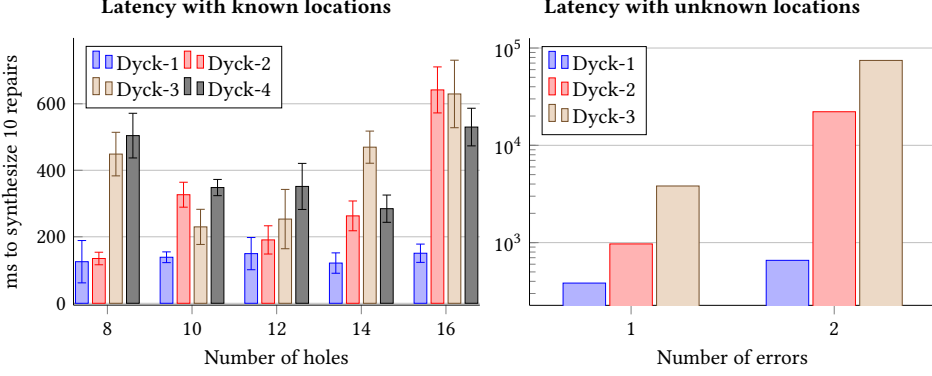


Dyck-1  $\rightarrow$  ( ) | ( Dyck-1 ) | Dyck-1 Dyck-1

Dyck-2  $\rightarrow$  Dyck-1 | [ ] | ( Dyck-2 ) | [ Dyck-2 ] | Dyck-2 Dyck-2

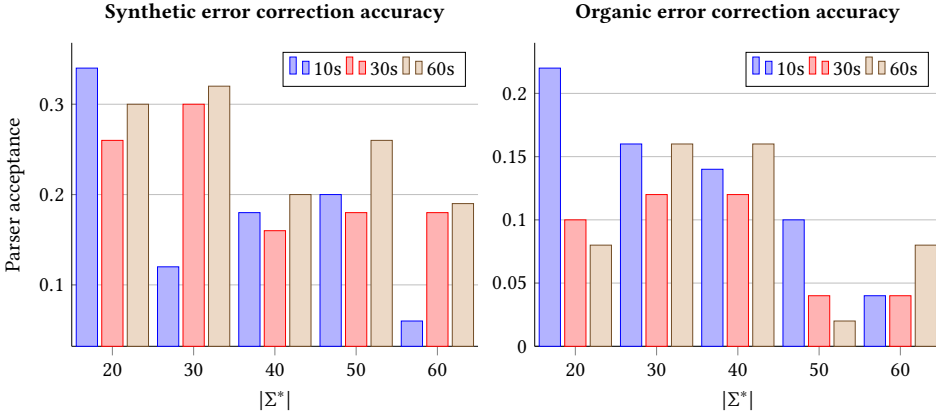
Dyck-3  $\rightarrow$  Dyck-2 | { } | ( Dyck-3 ) | [ Dyck-3 ] | { Dyck-3 } | Dyck-3 Dyck-3

In the first experiment, we sample a random valid string  $\sigma \sim \Sigma^{50} \cap \mathcal{L}_{\text{Dyck-n}}$ , then replace a fixed number of tokens with holes and measure the average time required to decode ten syntactically-admissible repairs across 100 trial runs.



In the second experiment, we sample a random valid string as before, but delete  $p$  tokens at random and rather than provide the location(s), ask our model to solve for both the location(s) and repair by sampling uniformly from all  $n$ -token HCs, then measure the total time required to decode the first admissible repair. Note the logarithmic scale on the y-axis.

In the following benchmark, we analyze error bracketing errors in a dataset of Java and Python code snippets mined from open-source repositories on GitHub. For Java, we sample valid single-line statements with bracket nesting more than two levels deep, synthetically delete one bracket uniformly at random, repair using Tidyparse<sup>3</sup>, then take the top-1 repair after  $t$  seconds, and validate using ANTLR’s Java 8 parser.



For Python, we sample invalid code fragments uniformly from the imbalanced bracket category of the Break-It-Fix-It (BIFI) dataset [33], a dataset of organic Python errors, repair using Tidyparse, take the top-1 repair after  $t$  seconds, and validate repairs using Python’s `ast.parse()` method.

## 9 DISCUSSION

While error correction with a few errors is tolerable, latency can vary depending on many factors including string length and grammar size. If errors are localized to the beginning or end of a string,

<sup>3</sup>Using the Dyck- $n$  grammar augmented with  $01 \rightarrow w \mid 01$ . Contiguous non-bracket tokens are substituted with a single placeholder token,  $w$ , and restored verbatim after bracket repair.

then latency is typically below 500ms. We observe that errors are typically concentrated nearby historical edit locations, which can be retrieved from the IDE or version control.

The primary advantage and limitation of our tool is the Levenshtein edit distance supported. This is an advantage, because in the absence of any contextual or statistical information, syntax and Levenshtein edits are often sufficiently constrained to determine a small number of valid repairs. This is also a limitation, because as the number of edits grows, the admissible set grows rapidly and the number of valid repairs may become too large to be useful depending on the language and snippet under repair.

Tidyparse in its current form has several other technical shortcomings: firstly it does not incorporate any neural language modeling technology at present, an omission we hope to address in the near future. Training a language model to predict likely repair locations and rank admissible results could lead to lower overall latency and more natural repairs. We also hope to explore the use of Metropolis-Hastings sampling to generate a diverse set of repairs, which could be useful for code completion and other applications.

Secondly, our current method generates sketch templates using a naïve enumerative search, feeding them individually to the SAT solver, which has the tendency to duplicate prior work and introduces unnecessary thrashing. Considering recent extensions of Boolean matrix-based parsing to linear context-free rewriting systems (LCFRS) [11], it may be feasible to search through these edits within the SAT solver, leading to yet unrealized and possibly significant speedups.

Lastly and perhaps most significantly, Tidyparse does not incorporate any semantic constraints, so its repairs whilst syntactically admissible, are not guaranteed to be semantically valid. This can be partly alleviated by filtering the results through an incremental compiler or linter, however the latency introduced may be non-negligible. It is also possible to encode type-based semantic constraints into the solver and intend to explore this direction more fully in future work.

We envision three primary use cases: (1) helping novice programmers become more quickly familiar with a new programming language (2) autocorrecting common typos among proficient but forgetful programmers and (3) as a prototyping tool for PL designers and educators. Featuring a grammar editor and built-in SAT solver, Tidyparse helps developers navigate the language design space, visualize syntax trees, debug parsing errors and quickly generate simple examples and counterexamples for testing.

## 10 CONCLUSION

The great compromise in program synthesis is one of efficiency versus expressiveness: the more expressive a language, the more concise and varied the programs it can represent, but the harder those programs are to synthesize without resorting to domain-specific heuristics. Likewise, the simpler a language is to synthesize, the weaker its concision and expressive power.

Most existing work on program synthesis focus on general  $\lambda$ -calculi, or very narrow languages such as finite automata or regular expressions. The former are too expressive to be synthesized or verified, whilst the latter are too restrictive to be useful. In our work, we focus on context-free and mildly context-sensitive grammars, which are expressive enough to capture a variety of useful programming language features, but not so expressive as to be unsynthesizable.

The second great compromise in program synthesis is that of reusability versus specialization. In programming, as in human affairs, there is a vast constellation of languages, each requiring specialized generators and interpreters. Are these languages truly irreconcilable? Or, as Noam Chomsky argues, are these merely dialects of a universal language? *Synthesis* then, might be a misnomer, and more aptly called *recognition*, in the analytic tradition.

In our work, we argue these two compromises are not mutually exclusive, but complementary and reciprocal. Programs and the languages they inhabit are indeed synthetic, but can be analyzed

and reused in the metalanguage of context-free grammars closed under conjunction. Not only does this admit an efficient synthesis algorithm, but allows users to introduce additional constraints without breaking compositionality, one of the most sacred tenets in programming language design.

Not only is linear algebra over finite fields an expressive language for inference, but also an efficient framework for inference on languages themselves. We illustrate a few of its applications for parsing incomplete strings and repairing syntax errors in context-free and sensitive languages. In contrast with LL and LR-style parsers, our technique can recover partial forests from invalid strings by examining the structure of  $M^*$  and handles arbitrary context-free languages. In future work, we hope to extend our method to more natural grammars like PCFG and LCFRS.

Furthermore, we argue it is possible to improve the efficiency of human programmers without sacrificing expressiveness by considering latency to synthesize an acceptable completion. In contrast with program synthesizers that require intermediate programs to be well-formed, our synthesizer is provably sound and complete up to a Levenshtein distance bound, and attempts to minimize total edits, but does not impose any constraints on the code itself being written.

Tidyparse accepts a CFG and a string to parse. If the string is valid, it returns the parse forest, otherwise, it returns a set of repairs, ordered by their Levenshtein edit distance to the invalid string. Our method compiles each CFG and candidate string onto a matrix dynamical system using an extended version of Valiant’s construction and solves for its fixedpoints using an incremental SAT solver. This approach to parsing has many advantages, enabling us to repair syntax errors, correct typos and generate parse trees for incomplete strings. By allowing the string to contain holes, repairs can contain either concrete tokens or nonterminals, which can be manually expanded by the user or a neural-guided search procedure. From a theoretical standpoint, this technique is particularly amenable to neural program synthesis and repair, naturally integrating with the masked-language-modeling task (MLM) used by transformer-based neural language models.

From a practical standpoint, we have implemented our approach as an IDE plugin and demonstrated its viability as a tool for live programming. Tidyparse is capable of generating repairs for invalid code in a range of practical languages. We plan to continue expanding its grammar and autocorrection functionality to cover a broader range of languages and hope to conduct a more thorough user study to validate its effectiveness in the near future.

## REFERENCES

- [1] Michael D Adams, Celeste Hollenbeck, and Matthew Might. 2016. On the complexity and performance of parsing with derivatives. In Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation. 224–236.
- [2] Alfred V Aho and Thomas G Peterson. 1972. A minimum distance error-correcting parser for context-free languages. SIAM J. Comput. 1, 4 (1972), 305–312.
- [3] Valentin Antimirov. 1996. Partial derivatives of regular expressions and finite automaton constructions. Theoretical Computer Science 155, 2 (1996), 291–319.
- [4] Roland Axelsson, Keijo Heljanko, and Martin Lange. 2008. Analyzing context-free grammars using an incremental SAT solver. In Automata, Languages and Programming: 35th International Colloquium, ICALP 2008, Reykjavik, Iceland, July 7-11, 2008, Proceedings, Part II 35. Springer, 410–422.
- [5] Yehoshua Bar-Hillel, Micha Perles, and Eli Shamir. 1961. On formal properties of simple phrase structure grammars. Sprachtypologie und Universalienforschung 14 (1961), 143–172.
- [6] Daniella Bar-Lev, Tuvi Etzion, and Eitan Yaakobi. 2021. On Levenshtein Balls with Radius One. In 2021 IEEE International Symposium on Information Theory (ISIT). 1979–1984. <https://doi.org/10.1109/ISIT45174.2021.9517922>
- [7] Leonor Becerra-Bonache, Colin de La Higuera, Jean-Christophe Janodet, and Frédéric Tantini. 2008. Learning Balls of Strings from Edit Corrections. Journal of Machine Learning Research 9, 8 (2008).
- [8] Janusz A Brzozowski. 1962. Canonical regular expressions and minimal state graphs for definite events. In Proc. Symposium of Mathematical Theory of Automata. 529–561.
- [9] Janusz A Brzozowski. 1964. Derivatives of regular expressions. Journal of the ACM (JACM) 11, 4 (1964), 481–494.
- [10] Janusz A. Brzozowski and Ernst Leiss. 1980. On equations for regular languages, finite automata, and sequential networks. Theoretical Computer Science 10, 1 (1980), 19–35.
- [11] Shay B Cohen and Daniel Gildea. 2016. Parsing linear context-free rewriting systems with fast matrix multiplication. Computational Linguistics 42, 3 (2016), 421–455.
- [12] Ryan Cotterell, Nanyun Peng, and Jason Eisner. 2014. Stochastic Contextual Edit Distance and Probabilistic FSTs. In Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, Vol. 2 (Short Papers). Association for Computational Linguistics, Baltimore, Maryland, 625–630.
- [13] Loris D’Antoni and Margus Veanes. 2014. Minimization of symbolic automata. In Proceedings of the 41st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages. 541–553.
- [14] Dingding Dong, Nitya Mani, and Yufei Zhao. 2023. On the number of error correcting codes. Combinatorics, Probability and Computing (2023), 1–14. <https://doi.org/10.1017/S0963548323000111>
- [15] Jay Earley. 1970. An efficient context-free parsing algorithm. Commun. ACM 13, 2 (1970), 94–102.
- [16] Seymour Ginsburg and H Gordon Rice. 1962. Two families of languages related to ALGOL. Journal of the ACM (JACM) 9, 3 (1962), 350–371.
- [17] E. T. Irons. 1963. An Error-Correcting Parse Algorithm. Commun. ACM 6, 11 (nov 1963), 669–673. <https://doi.org/10.1145/368310.368385>
- [18] Adam Kiezun, Vijay Ganesh, Philip J Guo, Pieter Hooimeijer, and Michael D Ernst. 2009. HAMPI: a solver for string constraints. In Proceedings of the eighteenth international symposium on Software testing and analysis. 105–116.
- [19] D. E. Knuth. 2005. Generating All Combinations and Partitions. Addison-Wesley, 5–6.
- [20] Lillian Lee. 2002. Fast context-free grammar parsing requires fast boolean matrix multiplication. Journal of the ACM (JACM) 49, 1 (2002), 1–15. <https://arxiv.org/pdf/cs/0112018.pdf>
- [21] Matthew Might, David Darais, and Daniel Spiewak. 2011. Parsing with derivatives: a functional pearl. ACM sigplan notices 46, 9 (2011), 189–195.
- [22] Alexander Okhotin. 2001. Conjunctive grammars. Journal of Automata, Languages and Combinatorics 6, 4 (2001), 519–535.
- [23] Rohit J. Parikh. 1966. On Context-Free Languages. J. ACM 13, 4 (oct 1966), 570–581. <https://doi.org/10.1145/321356.321364>
- [24] Perrin E. Ruth and Manuel E. Lladser. 2023. Levenshtein graphs: Resolvability, automorphisms determining sets. Discrete Mathematics 346, 5 (2023), 113310. <https://doi.org/10.1016/j.disc.2022.113310>
- [25] Itiroo Sakai. 1961. Syntax in universal translation. In Proceedings of the International Conference on Machine Translation and Applied Language Analysis.
- [26] Georgios Sakkas, Madeline Endres, Philip J Guo, Westley Weimer, and Ranjit Jhala. 2022. Seq2Parse: neurosymbolic parse error repair. Proceedings of the ACM on Programming Languages 6, OOPSLA2 (2022), 1180–1206.
- [27] Klaus U Schulz and Stoyan Mihov. 2002. Fast string correction with Levenshtein automata. International Journal on Document Analysis and Recognition 5 (2002), 67–85.
- [28] Matthew Szudzik. 2006. An elegant pairing function. In Special NKS 2006 Wolfram Science Conference. 1–12.



- [29] Michalis K Titsias and Christopher Yau. 2017. The Hamming ball sampler. J. Amer. Statist. Assoc. 112, 520 (2017), 1598–1611.
- [30] Leslie G Valiant. 1975. General context-free recognition in less than cubic time. Journal of computer and system sciences 10, 2 (1975), 308–315. <http://people.csail.mit.edu/virgi/6.s078/papers/valiant.pdf>
- [31] Leslie G Valiant. 1979. Completeness classes in algebra. In Proceedings of the eleventh annual ACM symposium on Theory of computing, 249–261. <https://dl.acm.org/doi/pdf/10.1145/800135.804419>
- [32] Alexander William Wong, Amir Salimi, Shaiful Chowdhury, and Abram Hindle. 2019. Syntax and Stack Overflow: A methodology for extracting a corpus of syntax errors and fixes. In 2019 IEEE International Conference on Software Maintenance and Evolution (ICSME). IEEE, 318–322.
- [33] Michihiro Yasunaga and Percy Liang. 2021. Break-it-fix-it: Unsupervised learning for program repair. In International Conference on Machine Learning. PMLR, 11941–11952.
- [34] Andreas Zeller. 2002. Isolating cause-effect chains from computer programs. ACM SIGSOFT Software Engineering Notes 27, 6 (2002), 1–10.
- [35] Qirun Zhang and Zhendong Su. 2017. Context-sensitive data-dependence analysis via linear conjunctive language reachability. In Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages. 344–358.

A EXAMPLE REPAIRS

We give some example human repairs that were correctly predicted from the StackOverflow dataset.

Original method	Human repair
<pre>from sympy import * x = Symbol('x', real=True) x, re(x), im(x)</pre>	<pre>from sympy import * x = Symbol('x', real=True) x, re(x), im(x)</pre>
<pre>result = yeald From(item.create()) raise Return(result)</pre>	<pre>result = yield From(item.create()) raise Return(result)</pre>
<pre>return 1/sum_p if sum_p \ return 0 else</pre>	<pre>return 1/sum_p if sum_p \ else 0</pre>
<pre>from itertools import permutations, product a='ABC' b=['*', '%3A'] l=[a]*2+[b] def g(list):     for p in permutations(list):         yield product(*p) result=g(l)</pre>	<pre>from itertools import permutations, product a='ABC' b=['*', '%3A'] l=[a]*2+[b] def g(list):     for p in permutations(list):         yield product(*p) result=g(l)</pre>
<pre>sum(len(v) for v items.values()))</pre>	<pre>sum(len(v) for v in items.values())</pre>
<pre>df.apply(lambda row: list(set(row['ids'])))</pre>	<pre>df.apply(lambda row: list(set(row['ids'])))</pre>
<pre>import numpy as np A_concat = np.array([a_0, a_1, a_2,..., a_n])</pre>	<pre>import numpy as np A_concat = np.array([a_0, a_1, a_2,..., a_n])</pre>