# 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. For example, we use it to repair syntax errors, perform sketch-based program synthesis, and decide various language induction and membership queries. This line of research provides an elegant unification of context-free program repair, code completion and sketch-based program synthesis. 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.

## 1 INTRODUCTION

Modern research on error correction 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, whether due to collision with a high-energy proton, manipulation by an adversary or some 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 code which ensures that even if some portion of the message should be corrupted through accidental or intentional means, one can still recover the original message by solving a linear system of equations. In particular, we frame our work inside the context of errors arising from human factors in computer programming.

In programming, most such errors initially manifest as syntax errors, and though often cosmetic, manual repair can present a significant challenge for novice programmers. The ECC problem may be refined by introducing a language, $\mathcal{L} \subset \Sigma^*$ and considering admissible edits transforming an arbitrary string, $s \in \Sigma^*$ into a string, $s' \in \mathcal{L}$. Known as *error-correcting parsing* (ECP), this problem was well-studied in the early parsing literature by Aho and Peterson [1], but fell out of favor for many years, perhaps due to its perceived complexity. By considering only minimal-length edits, ECP can be reduced to the so-called *language edit distance* (LED) problem, recently shown to be subcubic [3], suggesting its tractability. Previous LED results were primarily of a theoretical nature, but now, thanks to our contributions, we have finally realized a practical prototype.

## 2 TOY 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 §3, this is automatically 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 such as the one below, Tidyparse 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, Tidyparse 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
```

In the following paper, we will describe how we built it.

## 3  MATRIX THEORY

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

$$\mathcal{G} := \big\{ S \to S\,S \mid (\,S\,) \mid (\,) \big\} \Longrightarrow \mathcal{G}' = \big\{ S \to Q\,R \mid S\,S \mid L\,R, \quad L \to (, \quad R \to), \quad Q \to L\,S \big\}$$

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

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

If we define $\sigma_r^{\rightarrowtail} := \{w \mid (w \to \sigma_r) \in P\}$, then construct a matrix with nonterminals on the superdiagonal representing each token, $M^0_{r+1=c}(\mathcal{G}', e) := \sigma_r^{\rightarrowtail}$ and solve for the fixpoint $M^* = M + M^2$,
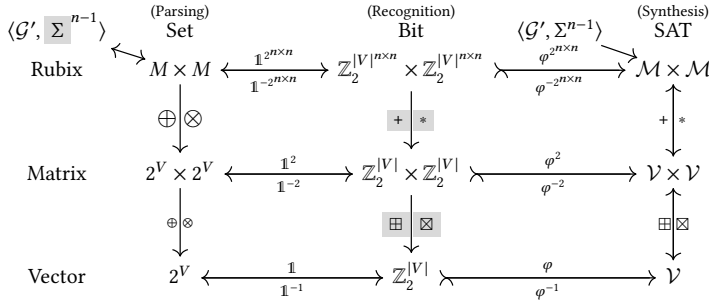
$$M^0 := \begin{pmatrix} \varnothing & \sigma_1^{\rightarrow} & \varnothing & \cdots & \varnothing \\ & & \ddots & & \vdots \\ & & & \ddots & \varnothing \\ & & & & \sigma_n^{\uparrow} \\ \varnothing & \cdots & \cdots & \cdots & \varnothing \end{pmatrix} \Rightarrow \ldots \Rightarrow M^* = \begin{pmatrix} \varnothing & \sigma_1^{\rightarrow} & \Lambda & \cdots & \Lambda_\sigma^* \\ & & \ddots & & \vdots \\ & & & \ddots & \Lambda \\ & & & & \sigma_n^{\uparrow} \\ \varnothing & \cdots & \cdots & \cdots & \varnothing \end{pmatrix}$$

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

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 $\varphi(\Lambda \circledast \Lambda') \equiv \varphi(\Lambda) \circledast \varphi(\Lambda')$ is suitable. This particular representation shares the same algebraic structure, but is more widely studied in coding theory, and easily translated into circuits for hardware acceleration. 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 [9], who first realized its time complexity was subcubic $\mathcal{O}(n^\omega)$ where $\omega$ is the matrix multiplication bound ($\omega < 2.77$), and Lee [6], who shows that speedups to CFL parsing are also realizable by Boolean matrix multiplication algorithms. Assuming sparsity, this technique is typically linearithmic, and is believed to be the most efficient general procedure for CFL recognition.

Valiant's decision procedure can be abstracted by lifting into the domain of bitvector variables, i.e., linear equations over finite fields, producing an algebraic expression for each nonterminal inhabitant of the northeasternmost bitvector $\mathcal{T}$, whose solutions correspond to valid parse forests for an incomplete string on the superdiagonal. This yields us a novel interpretation of Valiant's algorithm as an equational theory over finite fields, allowing us to solve for admissible hole inhabitants and parse forests. In particular, the following diagram commutes,[1]
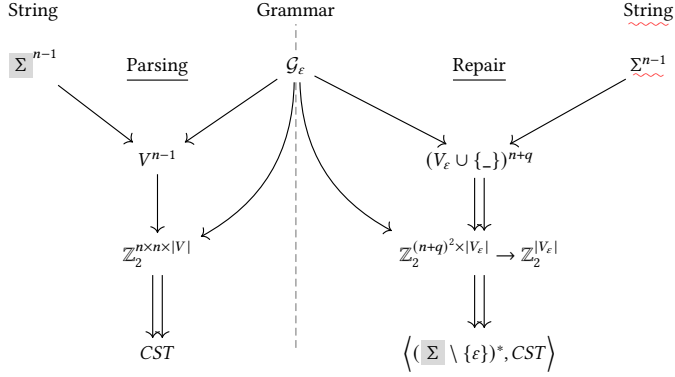


where $\mathcal{V}$ is a function $\mathbb{Z}_2^{|V|} \to \mathbb{Z}_2$. Note that while always possible to encode $\mathbb{Z}_2^{|V|} \to \mathcal{V}$ using the identity function, an arbitrary $\mathcal{V}$ might have zero, one, or in general, multiple solutions in $\mathbb{Z}_2^{|V|}$.

So far, we have only considered bounded context free languages with holes, however, we can extend this construction to the family of CFLs closed under conjunction, allowing us to incorporate more semantic properties and solve for variables in weakly context-sensitive languages. This extension will be covered in further depth in Sec. 11.

---

[1]Hereinafter, we use gray highlighting to denote types and functions defined over strings and binary constants only.
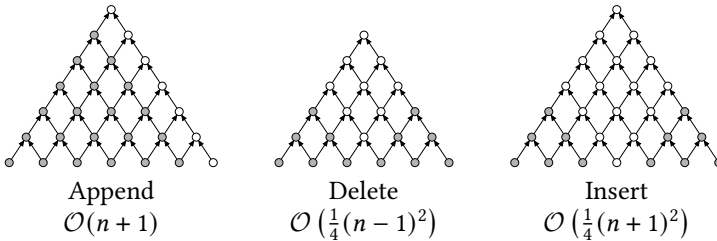
### 3.1 Encoding CFG parsing as SAT solving



Our algorithm produces set of concrete syntax trees (CSTs) for a given valid string. Otherwise, if the string contains an error, the algorithm generates a set of admissible corrections, alongside their CSTs.

$$M^* = \begin{pmatrix} \varnothing & \sigma_1^{\nrightarrow} & \mathcal{L}_{1,3} & \mathcal{L}_{1,3} & \mathcal{V}_{1,4} & \cdots & \mathcal{V}_{1,n} \\ & & \sigma_2^{\nrightarrow} & \mathcal{L}_{2,3} & & & \\ & & & \sigma_3^{\nrightarrow} & & & \\ & & & & \mathcal{V}_{4,4} & & \\ & & & & & \ddots & \\ & & & & & & \mathcal{V}_{n,n} \\ \varnothing & \cdots & & & & & \varnothing \end{pmatrix}$$

Depicted above is a SAT tensor representing $\sigma_1\ \sigma_2\ \sigma_3$ _ $\ldots$ _ 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 seeded nonterminal forests to be grown). Since $\mathcal{L}_{r,c}$ are fixed, we precompute them outside the SAT solver.

### 3.2 Incrementalizing the parser

Should only need to recompute submatrices affected by individual edits. In the worst case, each edit requires quadratic complexity in terms of $|\Sigma^*|$, assuming $\mathcal{O}(1)$ cost for each CNF-nonterminal subset join, $\mathbf{V}_1' \otimes \mathbf{V}_2'$.



| Append | Delete | Insert |
|---|---|---|
| $\mathcal{O}(n+1)$ | $\mathcal{O}\left(\frac{1}{4}(n-1)^2\right)$ | $\mathcal{O}\left(\frac{1}{4}(n+1)^2\right)$ |

Related to **dynamic matrix inverse** and **incremental transitive closure** with vertex updates. With a careful encoding, we can incrementally update SAT constraints as new keystrokes are received to eliminate redundancy.

## 4  ERROR RECOVERY

Unlike classical parsers which need special care to recover from errors, if the input string does not parse, Tidyparse can return partial subtrees. If no solution exists, the upper triangular entries will appear as a jagged-shaped ridge whose peaks represent the roots of parsable ASTs. These provide a natural debugging environment to aid the repair process.

The matrix $\mathcal{M}^*$ encodes a superposition of all labeled binary trees of a fixed breadth. Consider the string $\_ \; \dots \; \_$, which might generate various parse trees, with labels omitted to emphasize their structure:



Occasionally, it is not possible to decode a full tree. In a typical parser, error recovery requires special tricks. Here, we simply analyze the structure of the matrix $\mathbf{M}$ to decode the parseable subtrees:
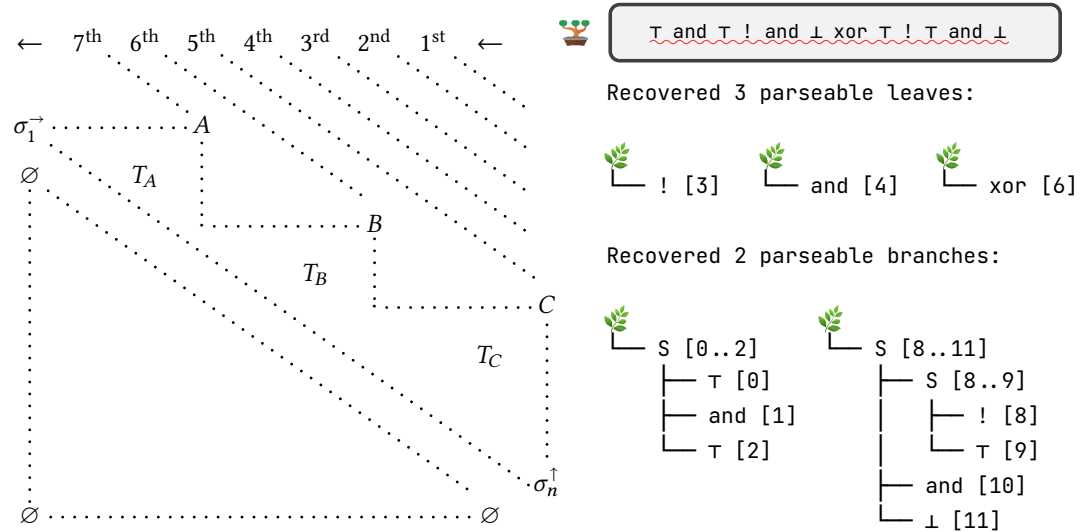


Fig. 1. Peaks along the UT matrix ridge correspond to maximally parseable substrings. By recursing over upper diagonals of decreasing elevation and discarding all subtrees that fall under the shadow of another's canopy, we can recover the partial subtrees. The example depicted above contains three such branches, rooted at nonterminals $C, B, A$.

These branches correspond to peaks on the upper triangular (UT) matrix ridge. As depicted in Fig. 1, we traverse the peaks by decreasing elevation to collect partial AST branches.

## 5 TREE DENORMALIZATION

Our parser emits a binary forest consisting of parse trees for the candidate string according to the CNF grammar, however this forest contains many so-called *Krummholz*, or *flag trees*, often found clinging to windy ridges and mountainsides.

**Algorithm 1** Tree denormalization

> **procedure** CUT(t: Tree)
>     stems ← {CUT(c) | c ∈ t.children}
>     **if** t.root ∈ $(V_{\mathcal{G}'} \setminus V_{\mathcal{G}})$ **then**
>         **return** stems
>     **else**
>         **return** { Tree(root, stems) }
>     **end if**
> **end procedure**



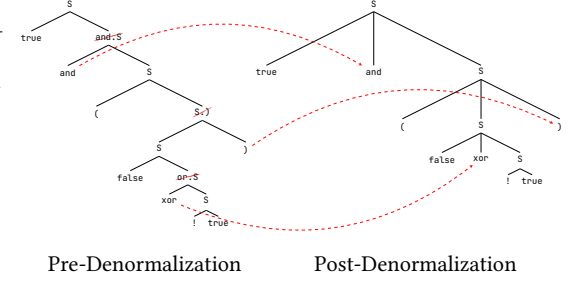Pre-Denormalization        Post-Denormalization

Fig. 2. 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).
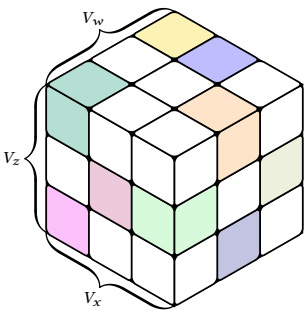
## 6 COARSENING AND SIMPLIFICATION OF SUBEXPRESSIONS

Clearly, the solving process is polynomial in the size of the string. For the sake of complexity, it would be convenient if well-formed subexpressions could be collapsed into a single nonterminal, or multiple nonterminals (in case of ambiguity). Naturally, this raises the question of when can partial derivations be simplified? This transformation is admissible and indeed desirable when the subexpression is "complete", i.e., its derivation cannot be altered by appending or prepending text. Returning to Fig. 4, under what circumstances is the following reduction admissible?

$$
\mathbf{M}^* = \begin{pmatrix} \varnothing & \sigma_1^{\pitchfork} & \mathcal{L}_{1,3} & \mathcal{L}_{1,3} & \mathcal{V}_{1,4} & \cdots & \mathcal{V}_{1,n} \\ & & \sigma_2^{\pitchfork} & \mathcal{L}_{2,3} & & & \\ & & & \sigma_3^{\pitchfork} & & & \\ & & & & \mathcal{V}_{4,4} & & \\ & & & & & \ddots & \\ & & & & & & \mathcal{V}_{n,n} \\ \varnothing & \cdots & & & & & \varnothing \end{pmatrix} \equiv \begin{pmatrix} \varnothing & \mathcal{L}_{1,3} & \mathcal{V}_{1,4} & \cdots & \mathcal{V}_{1,n} \\ & & \mathcal{V}_{2,4} & & \\ & & & \ddots & \\ & & & & \mathcal{V}_{n,n} \\ \varnothing & \cdots & & & \varnothing \end{pmatrix}
$$

For example, the string ( - b ) is morally *complete* in the sense that inserting adjacent text should not alter the interior derivation, while - 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. How can we formalize this argument to work on arbitrary (e.g., potentially ambiguous) grammars? This question can be reduced to a quotient: Does there exist another nonterminal, when joined with the substring in question, that will "strip away" any of its tokens, leading to another derivation? To answer this question, we must introduce the concept of a context-free language language derivative.

$$o \rightarrow \boxed{so} \mid \boxed{rs} \mid \boxed{rr} \mid \boxed{oo}$$
$$r \rightarrow \boxed{so} \mid \boxed{ss} \mid \boxed{rr} \mid \boxed{os}$$
$$s \rightarrow \boxed{so} \mid \boxed{rs} \mid \boxed{or} \mid \boxed{oo}$$

$$\mathcal{H}_{\{o\}} = \begin{pmatrix} \frac{\partial^2 o}{\partial \overleftarrow{o} \partial \overleftarrow{o}} & \frac{\partial^2 o}{\partial \overleftarrow{o} \partial \overrightarrow{r}} & \frac{\partial^2 o}{\partial \overleftarrow{o} \partial \overrightarrow{s}} \\ \frac{\partial^2 o}{\partial \overleftarrow{r} \partial \overrightarrow{o}} & \frac{\partial^2 o}{\partial \overleftarrow{r} \partial \overrightarrow{r}} & \frac{\partial^2 o}{\partial \overleftarrow{r} \partial \overrightarrow{s}} \\ \frac{\partial^2 o}{\partial \overleftarrow{s} \partial \overrightarrow{o}} & \frac{\partial^2 o}{\partial \overleftarrow{s} \partial \overrightarrow{r}} & \frac{\partial^2 o}{\partial \overleftarrow{s} \partial \overrightarrow{s}} \end{pmatrix}$$

$$\mathcal{H}_{\{r\}} = \begin{pmatrix} \frac{\partial^2 r}{\partial \overleftarrow{o} \partial \overleftarrow{o}} & \frac{\partial^2 r}{\partial \overleftarrow{o} \partial \overrightarrow{r}} & \frac{\partial^2 r}{\partial \overleftarrow{o} \partial \overrightarrow{s}} \\ \frac{\partial^2 r}{\partial \overleftarrow{r} \partial \overrightarrow{o}} & \frac{\partial^2 r}{\partial \overleftarrow{r} \partial \overrightarrow{r}} & \frac{\partial^2 r}{\partial \overleftarrow{r} \partial \overrightarrow{s}} \\ \frac{\partial^2 r}{\partial \overleftarrow{s} \partial \overrightarrow{o}} & \frac{\partial^2 r}{\partial \overleftarrow{s} \partial \overrightarrow{r}} & \frac{\partial^2 r}{\partial \overleftarrow{s} \partial \overrightarrow{s}} \end{pmatrix}$$

$$\mathcal{H}_{\{s\}} = \begin{pmatrix} \frac{\partial^2 s}{\partial \overleftarrow{o} \partial \overleftarrow{o}} & \frac{\partial^2 s}{\partial \overleftarrow{o} \partial \overrightarrow{r}} & \frac{\partial^2 s}{\partial \overleftarrow{o} \partial \overrightarrow{s}} \\ \frac{\partial^2 s}{\partial \overleftarrow{r} \partial \overrightarrow{o}} & \frac{\partial^2 s}{\partial \overleftarrow{r} \partial \overrightarrow{r}} & \frac{\partial^2 s}{\partial \overleftarrow{r} \partial \overrightarrow{s}} \\ \frac{\partial^2 s}{\partial \overleftarrow{s} \partial \overrightarrow{o}} & \frac{\partial^2 s}{\partial \overleftarrow{s} \partial \overrightarrow{r}} & \frac{\partial^2 s}{\partial \overleftarrow{s} \partial \overrightarrow{s}} \end{pmatrix}$$

Fig. 3. CFGs are witnessed by a rank-3 tensor, whose inhabitants indicate CNF productions. Gradients in this setting effectively condition the parse tensor M by constraining the superposition of admissible parse forests.

## 7   BACKPROPAGATION OF ERROR

Valiant's $\otimes$-operator, which yields the set of productions unifying known factors in a binary CFG, naturally implies the existence of a left- and right-quotient, which yield the set of nonterminals that may appear the right or left side of a known factor and its corresponding root. In other words, a known factor not only implicates subsequent expressions that can be derived from it, but also adjacent factors that may be composed with it to form a given derivation.

| Valiant's $\otimes$-operator | Left Quotient | Right Quotient |
|---|---|---|
| $x \otimes z := \{\, w \mid (w \rightarrow xz) \in P \,\}$ | $\frac{\partial w}{\partial \overleftarrow{x}} := \{\, z \mid (w \rightarrow xz) \in P \,\}$ | $\frac{\partial w}{\partial \overrightarrow{z}} := \{\, x \mid (w \rightarrow xz) \in P \,\}$ |

The left quotient coincides with the derivative operator first proposed by Brzozowski [4] and Antimirov [2] over regular languages, lifted into the context-free setting (our work). When the root and LHS are fixed, e.g., $\frac{\partial S}{\partial \overleftarrow{x}} : (\overleftarrow{V} \rightarrow S) \rightarrow \overrightarrow{V}$ returns the set of admissible nonterminals to the RHS. One may also consider a gradient operator, $\overleftarrow{\nabla} S : (\overleftarrow{V} \rightarrow S) \rightarrow \overrightarrow{V}$, which simultaneously tracks the partials with respect to a set of multiple LHS nonterminals produced by a fixed root.

If the root itself is unknown, we can define an operator, $\mathcal{H}_{\mathcal{W} \subseteq \mathcal{V}} : (\overleftarrow{V} \times \overrightarrow{V} \times \mathcal{W}) \rightarrow (\overleftarrow{V} \times \overrightarrow{V} \rightarrow \mathcal{W})$, which tracks second-order partial derivatives for all roots in $\mathcal{W}$. Unlike differential calculus on smooth manifolds, partials in this calculus do not necessarily commute depending on the CFG. By allowing the matrix $\mathcal{M}^*$ to contain bitvector variables representing holes in $\sigma$, we obtain a set of multilinear equations whose solutions exactly correspond to the set of admissible repairs and their corresponding parse forests. Specifically, the repairs coincide with holes in the superdiagonal $\mathcal{M}^*_{r+1=c}$, and the parse forests occur along upper-triangular entries $\mathcal{M}^*_{r+1<c}$. In the follow case, $\mathcal{M}^*$ is left-constrained, although the holes may (in general) appear anywhere in $\sigma$:

## 8  STRING REPAIR

Given some unparseable string, i.e., $\sigma_1 \dots \sigma_n : \Sigma^n \cap \mathcal{L}(\mathcal{G})^{\complement}$, where should we put holes to obtain a parseable $\tilde{\sigma} \in \mathcal{L}(\mathcal{G})$? To estimate the effect of perturbing $\sigma$ on $\Lambda_\sigma^*$, one can either (1) backpropagate $\nabla S$ across upper-triangular entries of $\mathcal{M}^*$, or (2) stochastically sample *minibatches* $\boldsymbol{\sigma} : \Sigma^{n \pm q} \sim \Delta_q(\sigma)$ from the Levenshtein q-ball centered on $\sigma$, i.e., the space of all edits with Levenshtein distance $\leq q$, loosely analogous to a finite difference approximation. In other words, we seek:

$$\bigcup_{\boldsymbol{\sigma} \sim \Delta_k(\sigma)} \{\tilde{\sigma} \in \boldsymbol{\sigma} \mid S \in \Lambda_{\tilde{\sigma}}^*\} \tag{2}$$

Let us consider (2). Suppose $U : \mathbb{Z}_2^{m \times m}$ is a matrix whose structure is shown in Eq. 3, wherein $C$ is a primitive polynomial over $\mathbb{Z}_2^m$ with coefficients $C_{1\dots m}$ and semiring operators $\oplus := \veebar, \otimes := \wedge$:

$$U^t V = \begin{pmatrix} C_1 & \cdots\cdots\cdots\cdots\cdots & C_m \\ \top & \circ & \cdots\cdots\cdots & \circ \\ \circ & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ \circ & \cdots & \circ & \top & \circ \end{pmatrix}^t \begin{pmatrix} V_1 \\ \vdots \\ \\ \\ V_m \end{pmatrix} \tag{3}$$

Since $C$ is primitive, the sequence $\mathbf{S} = (U^{0\dots 2^m - 1}V)$ must have *full periodicity*, i.e., for all $i, j \in [0, 2^m)$, $\mathbf{S}_i = \mathbf{S}_j \Rightarrow i = j$. To uniformly sample $\boldsymbol{\sigma} \sim \Sigma^d$ without replacement, we form an injection $\mathbb{Z}_2^m \to \Sigma^d$, cycle over $\mathbf{S}$, then discard samples which index any nonexistent element(s) of $\Sigma$. This method will reject $(1 - |\Sigma| 2^{-\lceil \log_2 |\Sigma| \rceil})^d$ samples overall, and requires $\mathcal{O}(1)$ per sample and $\mathcal{O}(|\Sigma|^d)$ to cover $\Sigma^d$.

Next, to admit deletion, we augment $P$ with $(\varepsilon^+ \to \varepsilon \mid \varepsilon^+ \varepsilon^+)$ and replace each production $(w \to t)$ with $(w \to t\, \varepsilon^+ \mid \varepsilon^+ t \mid t)$.

$$\frac{\Gamma \vdash \varepsilon \in \Sigma}{\Gamma \vdash (\varepsilon^+ \to \varepsilon \mid \varepsilon^+ \varepsilon^+) \in P} \; \varepsilon\text{-DUP} \qquad \frac{\Gamma \vdash (A \to B) \in P}{\Gamma \vdash (A \to B\, \varepsilon^+ \mid \varepsilon^+ B \mid B) \in P} \; \varepsilon^+\text{-INT}$$

Finally, to generate $\boldsymbol{\sigma} \sim \Delta_q(\sigma)$, we enumerate hole configurations $H(\sigma, i) = \sigma_{1\dots i-1}\, \_\, \_\, \sigma_{i+1\dots n}$ for each $i \in \binom{n}{d}$ and $d \in 1 \dots q$, then solve for $\mathcal{M}^*$. If $S \in \Lambda_\sigma^*$? has at least one solution, each edit in each $\tilde{\sigma} \in \boldsymbol{\sigma}$ will match exactly one of seven patterns:

$$\text{Deletion} = \begin{cases} \sigma_1 \; \dots \; \gamma_1\, \gamma_2 \; \dots \; \sigma_n & \gamma_{1,2} = \varepsilon \end{cases}$$

$$\text{Substitution} = \begin{cases} \sigma_1 \; \dots \; \gamma_1\, \gamma_2 \; \dots \; \sigma_n & \gamma_1 \neq \varepsilon \wedge \gamma_2 = \varepsilon \\ \sigma_1 \; \dots \; \gamma_1\, \gamma_2 \; \dots \; \sigma_n & \gamma_1 = \varepsilon \wedge \gamma_2 \neq \varepsilon \\ \sigma_1 \; \dots \; \gamma_1\, \gamma_2 \; \dots \; \sigma_n & \{\gamma_1, \gamma_2\} \cap \{\varepsilon, \sigma_i\} = \varnothing \end{cases}$$

$$\text{Insertion} = \begin{cases} \sigma_1 \; \dots \; \gamma_1\, \gamma_2 \; \dots \; \sigma_n & \gamma_1 = \sigma_i \wedge \gamma_2 \notin \{\varepsilon, \sigma_i\} \\ \sigma_1 \; \dots \; \gamma_1\, \gamma_2 \; \dots \; \sigma_n & \gamma_1 \notin \{\varepsilon, \sigma_i\} \wedge \gamma_2 = \sigma_i \\ \sigma_1 \; \dots \; \gamma_1\, \gamma_2 \; \dots \; \sigma_n & \gamma_{1,2} = \sigma_i \end{cases}$$

This approach is tractable for $n \lesssim 100, q \lesssim 3$, however more complex repairs require a more efficient gradient estimator. In the following section, we will give a more efficient construction for generating and accepting $\Delta_q(\sigma)$ that does not require enumerating hole configurations.

## 9   LEVENSHTEIN REACHABILITY

Levenshtein reachability is recognized by the nondeterministic infinite automaton (NIA) whose topology $\mathcal{L} = \vcenter{\hbox{\includegraphics}}$ can be factored into a product of (a) the monotone Chebyshev topology $\varkappa$, equipped with horizontal transitions accepting $\sigma_i$ and vertical transitions accepting Kleene stars, and (b) the monotone knight's topology $\vcenter{\hbox{\includegraphics}}$, equipped with transitions accepting $\sigma_{i+2}$. The structure of this space is representable as an acyclic NFA [8], populated by accept states within radius $k$ of $q_{n,0}$, or equivalently, a left-linear CFG whose productions finitely instantiate the transition dynamics:
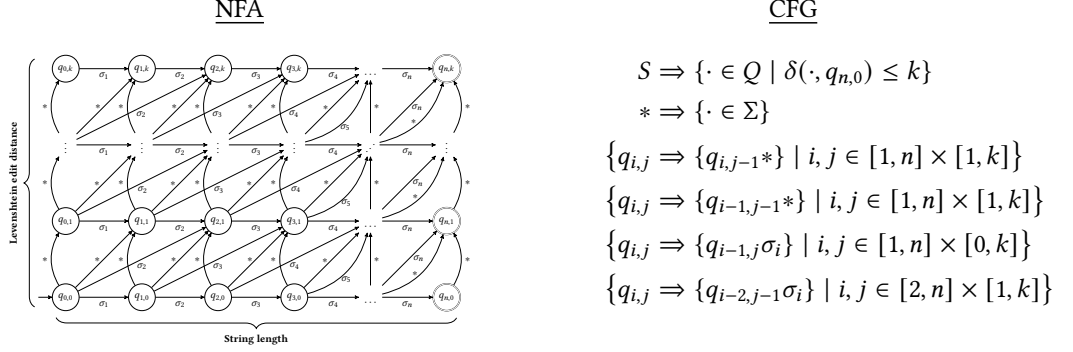
<u>NFA</u>



<u>CFG</u>

$$S \Rightarrow \{\cdot \in Q \mid \delta(\cdot, q_{n,0}) \leq k\}$$

$$* \Rightarrow \{\cdot \in \Sigma\}$$

$$\{q_{i,j} \Rightarrow \{q_{i,j-1}*\} \mid i, j \in [1, n] \times [1, k]\}$$

$$\{q_{i,j} \Rightarrow \{q_{i-1,j-1}*\} \mid i, j \in [1, n] \times [1, k]\}$$

$$\{q_{i,j} \Rightarrow \{q_{i-1,j}\sigma_i\} \mid i, j \in [1, n] \times [0, k]\}$$

$$\{q_{i,j} \Rightarrow \{q_{i-2,j-1}\sigma_i\} \mid i, j \in [2, n] \times [1, k]\}$$

Fig. 4. Levenshtein reachability from $\Sigma^n$ can be described as either an acyclic $\epsilon$-free NFA, or a left-linear CFG.

## 10   LANGUAGE EDIT REACHABILITY

Assume a hypothetical $\Phi(\mathcal{G}' : \mathbb{G}, \underset{\sim}{\sigma} : \Sigma^*) \mapsto \tilde{\sigma} : \mathcal{L}(\mathcal{G}')$ which takes a CFG, $\mathcal{G}'$, generating an arbitrary nonempty CFL, and an unparseable string, $\underset{\sim}{\sigma}$, and which returns element(s) of $\mathcal{L}(\mathcal{G}')$ most similar to $\underset{\sim}{\sigma}$ according to their Levenshtein distance $\Delta(\underset{\sim}{\sigma}, \cdot)$.

Let $G(\underset{\sim}{\sigma} : \Sigma^*, d : \mathbb{N}^+) \mapsto \mathbb{G}$ be the specific construction described in Sec. 9 which accepts a string, $\underset{\sim}{\sigma}$, and an edit distance, $d$, and returns a grammar representing the NFA that recognizes the language of all strings within Levenshtein radius $d$ of $\underset{\sim}{\sigma}$. To find the language edit distance and corresponding least-distance edit, we must find the smallest $d$ such that $\mathcal{L}_d^\cap$ is nonempty, where $\mathcal{L}_d^\cap$ is defined as $\mathcal{L}(G(\underset{\sim}{\sigma}, d)) \cap \mathcal{L}(\mathcal{G}')$. In other words, we seek $\tilde{\sigma}$ and $d^*$ under which three criteria are equisatisfiable: (1) $\tilde{\sigma} \in \mathcal{L}(\mathcal{G}')$, and (2) $\Delta(\underset{\sim}{\sigma}, \tilde{\sigma}) \leq d^* \iff \tilde{\sigma} \in \mathcal{L}(G(\underset{\sim}{\sigma}, d^*))$, and



Fig. 5. LED is computed gradually by incrementing $d$ until $\mathcal{L}_d^\cap \neq \varnothing$.

(3) $\nexists \sigma' \in \mathcal{L}(\mathcal{G}').[\Delta(\underset{\sim}{\sigma}, \sigma') < d^*]$. To satisfy these criteria, it suffices to check $d \in (1, d_{\max}]$ 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) $d_{\max}$ is attained. This procedure is guaranteed to terminate in at most either (1) the number of steps required to overwrite every symbol in $\underset{\sim}{\sigma}$, or (2) the length of the shortest string in $\mathcal{L}(\mathcal{G}')$, whichever is greater. More precisely:

$$\varphi_{d+1}(\mathcal{G}', \underset{\sim}{\sigma}) := \begin{cases} \varphi[\tilde{\sigma} \in \mathcal{L}(G(\underset{\sim}{\sigma}, d)) \wedge \tilde{\sigma} \in \mathcal{L}(\mathcal{G}')] & \text{if } d = 1 \text{ or SAT.} \\ \varphi_d \oplus \bigoplus_{\{q \in Q \mid \delta(q, q_{n,0}) = d+1\}} \varphi[S \to q] & \text{if } d \leq \max(|\underset{\sim}{\sigma}|, \min_{\sigma \in \mathcal{L}(\mathcal{G}')} |\sigma|). \end{cases} \tag{4}$$

The function $\varphi_{d+1}(\mathcal{G}', \underset{\sim}{\sigma})$ is a realizer of $\Phi$.
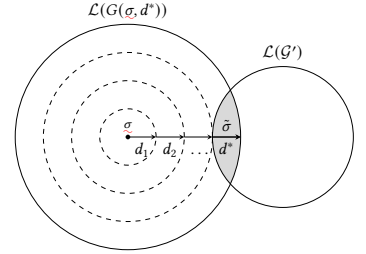
## 11  LINEAR CONJUNCTIVE REACHABILITY

It is well-known that the family of CFLs is not closed under intersection. For example, consider $\mathcal{L}_\cap := \mathcal{L}(\mathcal{G}_1) \cap \mathcal{L}(\mathcal{G}_2)$:

$$P_1 := \big\{\ S \to LR, \quad L \to ab \mid aLb, \quad R \to c \mid cR\ \big\}$$
$$P_2 := \big\{\ S \to LR, \quad R \to bc \mid bRc, \quad L \to a \mid aL\ \big\}$$

Note that $\mathcal{L}_\cap$ generates the language $\big\{\ a^d b^d c^d \mid d > 0\ \big\}$, which according to the pumping lemma is not context-free. We can encode $\bigcap_{i=1}^c \mathcal{L}(\mathcal{G}_i)$ 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 $E(t_\cap, \mathcal{G}_i) = \{w_i \mid (w_i \to t_\cap) \in P_i\}$ and bind them together:

$$\bigwedge_{t \in \Sigma_\cap} \bigwedge_{j=1}^{c-1} \bigwedge_{i=1}^{|\sigma|} E(t_\cap, \mathcal{G}_j) \equiv_{\sigma_i} E(t_\cap, \mathcal{G}_{j+1}) \tag{5}$$



Fig. 6. Orientations of a $\bigcap_{i=1}^4 \mathcal{L}(\mathcal{G}_i) \cap \Sigma^6$ configuration. As $c \to \infty$, this shape approximates a circular cone whose symmetric axis joins $\sigma_i$ with orthonormal unit productions $w_i \to t_\cap$, 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.

It is a well-known fact in formal language theory that CFLs are closed under union, composition, substitution and intersection with regular languages, and the intersection between two CFLs is decidable. However CFLs are not closed under intersection, which requires a more general formalism. Following Okhotin [7], we define higher-order grammar combinators $\cup, \cap : \mathcal{G}^* \times \mathcal{G}^* \to \mathcal{G}^*$ where $\mathcal{G}^*$ is a conjunctive grammar, which are basically finite collections of CFGs. Unlike parser combinators which are susceptible to ambiguity errors, our grammar combinators return parse forests in case of syntactic ambiguity, and do not suffer from the same shortcomings.

Given two CFLs $\mathcal{L}_\mathcal{G}, \mathcal{L}_{\mathcal{G}'}$, we can compute the intersection $\mathcal{L}_\mathcal{G} \cap \mathcal{L}_{\mathcal{G}'} \cap \Sigma^d$ by encoding $(\mathbf{M}_\mathcal{G}^* \sigma) = (\mathbf{M}_{\mathcal{G}'}^*, \sigma)$. This allows us to build a DSL of grammar combinators to constrain the solution space.

For example, we can solve $\Sigma^d \cap \overline{\mathcal{L}_\mathcal{G}}$ by enumerating $\{\beta\sigma'\gamma \mid \sigma' \in \Sigma^d, \beta = \gamma = \_^k\}$, overapproximating the prefix and suffix (padding left and right), and checking for UNSAT to underapproximate *impossible substrings*, strings which cannot appear in any $\{\sigma \in \mathcal{L}_\mathcal{G}\}$. Precomputing impossible substrings for a given grammar allows us to quickly eliminate inadmissible repairs and localize syntax errors in candidate strings.

Using the technique from Sec. 9, we can also compute language edit distance, the minimum number of Levenshtein edits required to fix a syntactically invalid string. Language intersection is significantly faster than approximating the gradient via sampling.

We can also build a set of grammars of increasing granularity, like a lattice structure. Basically, we can build up a lattice (in the order theoretic sense), consisting of grammars of increasing granularity. All programming languages require balanced parentheses, but some have additional constraints. So we can combine grammars, count and do bounded linear integer arithmetic.

## 12  PRACTICAL EXAMPLE

Tidyparse requires a grammar – this can be either provided by the user or ingested from a BNF-like specification. The following is a slightly more complex grammar, designed to resemble a more realistic use case:

```
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 VO Vexp
Vexp -> ( VarName , VarName ) | Vexp Vexp
VarName -> a | b | c | d | e | ... | z
FunName -> foldright | map | filter
VO -> + | - | * | / | > | = | < | `||` | &&
---
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> ) )
...
```

### 12.1  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>
...
```

### 12.2  Interactive Nonterminal Expansion

Users can interactively build up a complex expression by placing the caret over a placeholder they wish to expand,

```
if <Vexp> X then <Vexp> else <Vexp>
```

then invoking Tidyparse by pressing `ctrl` + `Space`, to receive a list of expressions constistient with the grammar:

```
if map X then <Vexp> else <Vexp>
if uncurry X then <Vexp> else <Vexp>
if foldright X then <Vexp> else <Vexp>
...
```
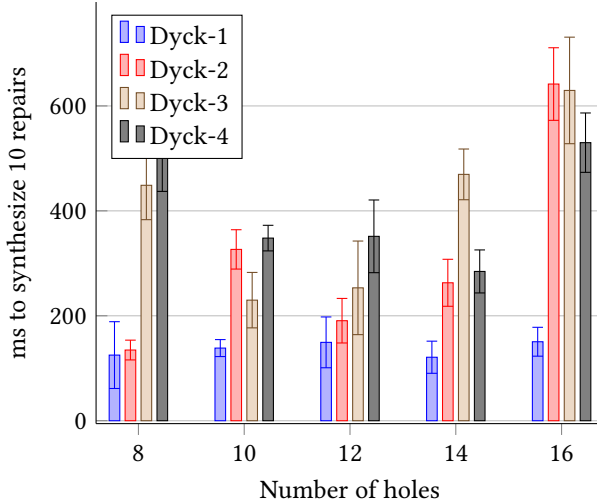
## 13  EVALUATION

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 types of balanced parentheses.
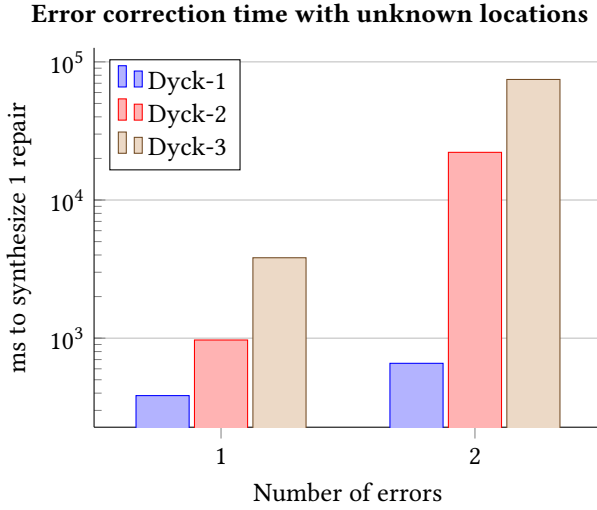
```
D1 -> () | ( D1 ) | D1 D1
D2 -> D1 | [ ] | ( D2 ) | [ D2 ] | D2 D2
D3 -> D2 | { } | ( D3 ) | [ D3 ] | { D3 } | D3 D3
```

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 tokens with holes and measure the average time taken to decode ten syntactically-admissible repairs across 100 trial runs.

**Error correction time with known locations**



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 the logarithmic scale on the y-axis.

**Error correction time with unknown locations**



## 13.1 Task Difficulty

**Level ½** - Incomplete code (all missing tokens at end)
**Level 1** - Known number of errors at known locations
**Level 2** - Known number of errors at unknown locations
**Level 3** - Unknown number of errors at unknown locations

## 13.2 Repair Datasets

Synthetic - Synthetic errors in natural code snippets
Natural - Natural errors and fixes mined from Git history

## 13.3 Evaluation Metrics

Syntactic validity - Only considers whether the repair parses
Repair alignment - Measures string distance to true fix

## 13.4 Model configurations

Model configurations - LLM, ECP, ECP + LLM GraphCodeBERT RobertA Incoder CarperAI?

## 14 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 known to be concentrated in specific locations, such as the beginning or end of a string, then latency is typically below 500ms. Should errors occur uniformly at random, admissible repairs can take longer to discover, however these scenarios are unusual in our experience. We observe that errors are typically concentrated nearby historical edit locations, which can be retrieved from the IDE or version control. Further optimizations that reduce the total number of repairs checked are possible by eliminating improbable sketch templates.

Tidyparse in its current form has a number of 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.

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) [5], 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 while syntactically admissible, are not guaranteed to be semantically valid. We note however, that it is possible to encode type-based semantic constraints into the solver and intend to explore this direction more fully in future work.

Although not intended to be a dedicated parser and we make no attempt to rigorously compare parsing latency, parsing valid strings with Tidyparse is typically competitive with classical parsing methods. Our primary motivation is to facilitate the usability and explainability of parsing with errors. 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 educators and designers.

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. Although the algorithm may seem esoteric at first glance, in our experience it is much more interpretable than classical parsers, which exhibit poor error-recovery and diagnostics.

## 15 CONCLUSION

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 toy 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] 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.

[2] Valentin Antimirov. 1996. Partial derivatives of regular expressions and finite automaton constructions. Theoretical Computer Science 155, 2 (1996), 291–319.

[3] Karl Bringmann, Fabrizio Grandoni, Barna Saha, and Virginia Vassilevska Williams. 2019. Truly subcubic algorithms for language edit distance and RNA folding via fast bounded-difference min-plus product. SIAM J. Comput. 48, 2 (2019), 481–512.

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

[5] Shay B Cohen and Daniel Gildea. 2016. Parsing linear context-free rewriting systems with fast matrix multiplication. Computational Linguistics 42, 3 (2016), 421–455.

[6] 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

[7] Alexander Okhotin. 2001. Conjunctive grammars. Journal of Automata, Languages and Combinatorics 6, 4 (2001), 519–535.

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

[9] 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

## A    EXAMPLE REPAIRS

| 1.a) Original method | 1.b) Synonymous variant |
|---|---|
| ```<br>public void flush(int b) {<br>    buffer.write((byte) b);<br>    buffer.compact();<br>}<br>``` | ```<br>public void flush(int b) {<br>    cushion.write((byte) b);<br>    cushion.compact();<br>}<br>``` |
| 2.a) Multi-masked method | 2.b) Multi-masked variant |
| ```<br>public void <MASK>(int b) {<br>    buffer.<MASK>((byte) b);<br>    <MASK>.compact();<br>}<br>``` | ```<br>public void <MASK>(int b) {<br>    cushion.<MASK>((byte) b);<br>    <MASK>.compact();<br>}<br>``` |
| 3.a) Model predictions | 3.b) Model predictions |
| ```<br>public void output(int b) {<br>    buffer.write((byte) b);<br>    buffer.compact();<br>}<br>``` | ```<br>public void append(int b) {<br>    cushion.add((byte) b);<br>    cushion.compact();<br>}<br>``` |