# Tidyparse: Real-Time Context Free Error Correction

Breandan Mark Considine
McGill University
bre@ndan.co

Jin Guo
McGill University
jguo@cs.mcgill.ca

Xujie Si
McGill University
xsi@cs.mcgill.ca

## Abstract

Tidyparse is a program synthesizer that performs real-time error correction for context free languages. Given both an arbitrary context free grammar (CFG) and an invalid string, the tool lazily generates admissible repairs while the author is typing, ranked by Levenshtein edit distance. Repairs are guaranteed to be sound, complete, syntactically valid and minimal. Tidyparse is the first system of its kind offering these guarantees in a real-time editor. 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. We have released an IDE plugin demonsrating the system described.[1]

## 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, cf. 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 [2], suggesting its possible tractability. Previous results on ECP and LED were primarily of a theorietical nature, but now, thanks to our contributions, we have finally realized a practical prototype.
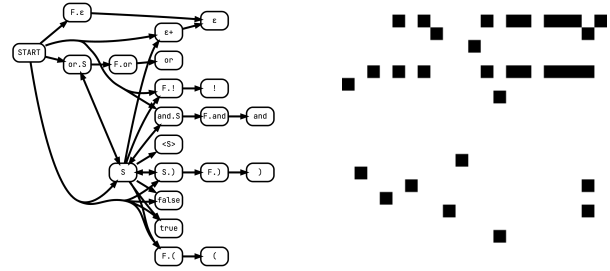
---

## 2  Toy Example

Suppose we are given the following context free grammar:

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

This gets rewritten into the following CNF grammar:

```
F.! → !    ε+ → ε       S → false     F.and → and
F.( → (    ε+ → ε+ ε+    S → F.! S       S.) → S F.)
F.) → )    S → <S>       S → S or.S    or.S → F.or S
F.ε → ε    S → true      S → S and.S   and.S → F.and S
F.or → or  S → S ε+       S → F.( S.)
```

We can visualize the CFG as either a graph or a matrix:
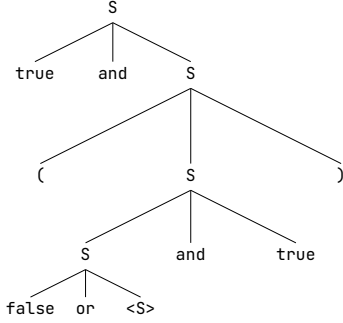


This CFG corresponds to the Dyck-3 language of balanced parentheses. Suppose we have the following string:

```
true and ( false or and true false
```

Our tool produces the following suggestions, where green is insertion, orange is substitution and red is deletion.

```
1.) true and ( false or ! true )
2.) true and ( false or <S> and true )
3.) true and ( false or ( true ) )
4.) true and ( false or ! ! false )
5.) true and ( false or <S> or false )
6.) true and ( false or ! <S> and true )
7.) true and ( false or ! ( false ) )
8.) true and ( false or ! true ) or <S>
9.) true and ( false or ! <S> ) and true false
```

Selecting the third entry will produce a single parse tree:

```
                    S
          ┌─────────┼─────────┐
        true       and        S
                      ┌────────┼────────┐
                      (        S        )
                        ┌───────┼───────┐
                        S      and     true
                   ┌────┼────┐
                 false  or   <S>
```

All the above visualizations were generated automatically.

## 3 Theory

It is a well-known fact that every CFG can be reduced into Chomsky Normal Form (CNF). Given a CFG, $\mathcal{G} : \langle V, \Sigma, P, S \rangle$ in CNF, we can construct a recognizer $R_{\mathcal{G}} : \Sigma^n \to \mathbb{B}$ for strings $\sigma : \Sigma^n$ as follows. Let $\mathcal{P}(V)$ be our domain, 0 be $\varnothing$, $\oplus$ be $\cup$, and $\otimes$ be defined as follows:

$$a \otimes b := \{C \mid \langle A, B \rangle \in a \times b, (C \to AB) \in P\}$$

We initialize $\mathbf{M}^0_{r,c}(\mathcal{G}, \sigma) := \{V \mid c = r + 1, (V \to \sigma_r) \in P\}$ and search for a matrix $\mathbf{M}^*$ via fixpoint iteration,

$$\mathbf{M}^* = \begin{pmatrix} \varnothing & \{V\}_{\sigma_1} & \dots & \dots & \mathcal{T} \\ \varnothing & \varnothing & \{V\}_{\sigma_2} & \dots & \dots \\ \varnothing & \varnothing & \varnothing & \{V\}_{\sigma_3} & \dots \\ \varnothing & \varnothing & \varnothing & \varnothing & \{V\}_{\sigma_4} \\ \varnothing & \varnothing & \varnothing & \varnothing & \varnothing \end{pmatrix}$$

where $\mathbf{M}^*$ is the least solution to $\mathbf{M} = \mathbf{M} + \mathbf{M}^2$. We can then define the recognizer as $S \in \mathcal{T}? \iff \sigma \in \mathcal{L}(\mathcal{G})$?
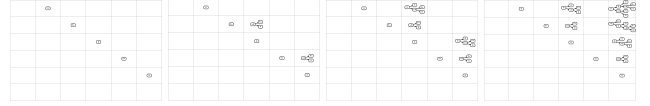
While theoretically elegant, this decision procedure can be optimized by lowering onto a rank-3 binary tensor. We do so simply by noting that $\bigoplus_{k=1}^n \mathbf{M}_{ik} \otimes \mathbf{M}_{kj}$ has cardinality bounded by $|V|$ and is thus representable as a fixed-length vector using the characteristic function, $\mathbb{1}$. In particular, $\oplus, \otimes$ are defined as $\boxplus, \boxtimes$, so that the following diagram commutes:

$$\begin{array}{ccc} V \times V & \xrightarrow{\oplus/\otimes} & V \\ \mathbb{1}^{-2} \big\Uparrow \mathbb{1}^2 & & \mathbb{1}^{-1} \big\Uparrow \mathbb{1} \\ \mathbb{B}^{|V|} \times \mathbb{B}^{|V|} & \xrightarrow[\boxplus/\boxtimes]{} & \mathbb{B}^{|V|} \end{array}$$

Full details of this bisimilarity can be found in Valiant [4] and Lee [3], who proves its time complexity to be $\mathcal{O}(n^\omega)$ where $\omega$ is the matrix multiplication bound. By assuming sparsity, this technique can typically be reduced to linearithmic time, and is currently the best known asymptotic bound for context free language recognition to date.

By allowing the matrix $\mathbf{M}^0_{r,c}$ to contain bitvector variables $\mathcal{B}^{|V|}$ representing holes in the string, we obtain a set of multilinear equations whose solutions exactly correspond to the set of admissible repairs and their corresponding parse trees.

Specifically, the repairs occur along holes in the superdiagonal $\mathbf{M}^*_{r+1=c}$, and the upper-triangular entries $\mathbf{M}^*_{r+1<c}$ represent the corresponding parse tree. If no solution exists, then the upper triangular entries will appear as a jagged-shaped ridge whose peaks represent the roots of the parsable subtree. We illustrate this fact in the following example:



## 4 Implementation

Tidyparse accepts a CFG and a string to parse and returns a set of candidate strings, ordered by their Levenshtein edit distance to the original string. Our method lowers the CFG and candidate string onto a matrix dynamical system using an extended version of Valiant's construction and solves for the fixpoint matrix using an incremental SAT solver.

Here is how we implemented it (and so can you) [4].

## 5 Examples

There are examples too.

## 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] 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.

[3] 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

[4] 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