

Propagation of Syntax Errors in Context-Sensitive Languages using the Lifted Brzozowski Derivative

Breandan Mark Considine
McGill University
bre@ndan.co

Xujie Si
McGill University
xsi@cs.mcgill.ca

Abstract

Brzozowski defines the derivative of a regular language as the suffixes that complete a known prefix. In this work, we establish a Galois connection between Brzozowski’s derivative and Valiant’s fixpoint construction in the context-free setting, and further extend their work into the hierarchy of bounded context-sensitive languages realizable by finite CFL intersection. We show how context-sensitive language recognition can be reduced into a tensor algebra over finite fields, drawing a loose analogy to partial differentiation in Euclidean spaces. In addition to its theoretical contributions, our method has yielded applications to incremental parsing, code completion and program repair. For example, we use it to repair syntax errors and perform sketch-based program synthesis, among other language decision problems.

1 Introduction

Recall that a CFG is a quadruple consisting of terminals (Σ), 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 $w \rightarrow xy$, or $w \rightarrow \sigma$, where $w, x, y: V$ and $\sigma: \Sigma$. For example, the CFG, $P = \{S \rightarrow SS \mid (S) \mid ()\}$, corresponds to the CNF:

$$P' = \{ S \rightarrow XR \mid SS \mid LR, \quad L \rightarrow (, \quad R \rightarrow), \quad X \rightarrow LS \}$$

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

$$x \otimes y := \{ W \mid \langle X, Y \rangle \in x \times y, (W \rightarrow XY) \in P \} \quad (1)$$

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

$$\mathbf{M}^* = \begin{pmatrix} \emptyset & \{V\}_{\sigma_1} & \dots & \mathcal{T} \\ \vdots & \vdots & \ddots & \vdots \\ \emptyset & \dots & \{V\}_{\sigma_n} & \emptyset \end{pmatrix} \quad (2)$$

where \mathbf{M}^* is the least solution to $\mathbf{M} = \mathbf{M} + \mathbf{M}^2$. We can then define the recognizer as: $S \in \mathcal{T}? \iff \sigma \in \mathcal{L}(\mathcal{G})?$

Full details of the bisimilarity between parsing and matrix multiplication can be found in Valiant [4] and Lee [3], who shows its time complexity to be $\mathcal{O}(n^\omega)$ where ω is the least matrix multiplication upper bound (currently, $\omega < 2.77$).

Note 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:

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

This construction can be lifted into the domain of bitvector variables, producing an algebraic expression for each scalar inhabitant of the northeasternmost bitvector, whose solutions correspond to valid parse forests for an incomplete string on the superdiagonal. Consider two cases, where the derivation is left- or right-constrained, i.e., $\alpha \gamma, \gamma \alpha$.¹

Valiant’s \otimes operator, which solves for the set of productions unifying known factors in a binary CFG, implies the existence of a left- and right-quotient, which yield the set of nonterminals that may appear to the right- or left-side, respectively, of known factors in a ternary configuration.

Left Quotient

Right Quotient

$$\frac{\partial f}{\partial x} = \{ y \mid (w \rightarrow xy) \in P \} \quad \frac{\partial f}{\partial y} = \{ x \mid (w \rightarrow xy) \in P \}$$

x	w
	y

x	w
	y

The left quotient coincides with the derivative operator first proposed by Brzozowski [2] and Antimirov [1] over regular languages, lifted into the context-free setting (our work).

Let \mathcal{V} represent 2^V . Assuming the root is known, (e.g., S), the operator $\frac{\partial S}{\partial x} : (\tilde{V} \rightarrow S) \rightarrow \tilde{V}$ is a dependently-typed function which returns the unknown RHS factor. We may also define a gradient operator, $\tilde{\nabla} S : (\tilde{V} \rightarrow S) \rightarrow \tilde{V}$, which tracks the partials with respect to multiple unknown factors.

If the root itself is unknown, we can define an operator, $\mathcal{H}_{\mathcal{W} \subseteq \mathcal{V}} : (\tilde{V} \times \tilde{V} \times \mathcal{W}) \rightarrow (\tilde{V} \times \tilde{V} \rightarrow \mathcal{W})$, 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.

¹Hereinafter, we shall use gray highlighting to distinguish between bound variables (i.e., constants) and free variables that are unhighlighted.

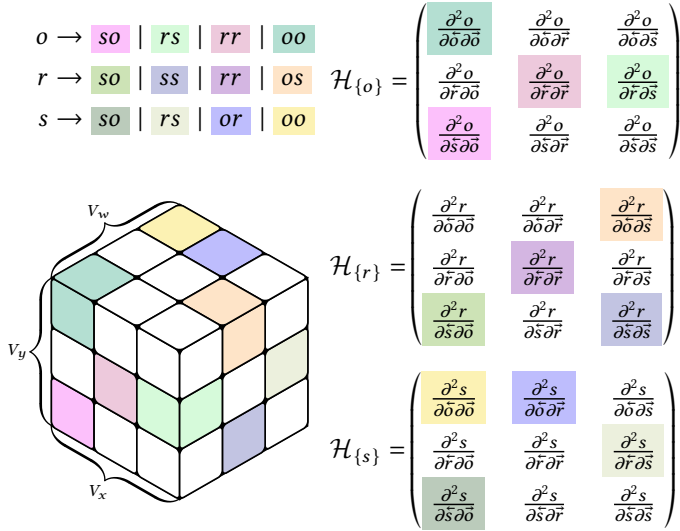


Figure 1. CFGs are witnessed by a Rank-3 binary tensor, whose nonzero entries indicate CNF productions. Derivatives in this setting effectively condition the parse tensor. By backpropagating \mathcal{H} across upper-triangular entries of \mathbf{M}^* , we constrain the superposition of admissible parse forests.

2 Context Sensitive Reachability

A well-known result in formal language theory is that the family of CFLs is not closed under intersection. For example, let us consider $\mathcal{L}^\cap := \mathcal{L}(\mathcal{G}_1) \cap \mathcal{L}(\mathcal{G}_2)$, defined as follows:

$$\begin{aligned} \mathcal{G}_1 &:= \{ S \rightarrow LR \quad L \rightarrow ab \mid aLb \quad R \rightarrow c \mid cR \} \\ \mathcal{G}_2 &:= \{ S \rightarrow LR \quad R \rightarrow bc \mid bRc \quad L \rightarrow a \mid aL \} \end{aligned}$$

Note that \mathcal{L}^\cap generates the language $\{ a^n b^n c^n \mid n > 0 \}$, which according to the pumping lemma is not context free. In our formalism, we encode bounded intersections of CFLs $\bigcap_{\mathcal{G} \in \Gamma} \mathcal{L}(\mathcal{G})$ as an elongated polygonal prism with upper-triangular matrices adjoined to each rectangular face, like the fletching on an arrow. The center shaft represents equivalent nonterminals on the left hand side of a unit production. Solutions to this system of equations yield bounded-width strings at the intersection of a finite collection of CFLs.

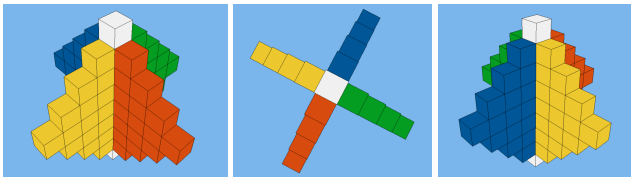


Figure 2. Orientations of a $\bigcap_{i=1}^4 \mathcal{L}(\mathcal{G}_i) \cap \Sigma^9$ configuration.

2.1 Encoding CFG parsing as SAT solving

By allowing the matrix \mathbf{M}^* in Eq. 2 to contain bitvector variables representing holes in the string and nonterminal sets,

we obtain a set of multilinear SAT 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 $\mathbf{M}_{r+1=c}^*$, and the parse forests occur along the upper-triangular entries $\mathbf{M}_{r+1 < c}^*$.

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

2.2 Deletion, substitution, and insertion

Deletion, substitution and insertion can be simulated by adding a left- and right- ε production to each unit production:

$$\begin{aligned} &\frac{\Gamma \vdash \varepsilon \in \Sigma}{\Gamma \vdash (\varepsilon^+ \rightarrow \varepsilon \mid \varepsilon^+ \varepsilon^+) \in P} \varepsilon\text{-DUP} \\ &\frac{\Gamma \vdash (A \rightarrow B) \in P}{\Gamma \vdash (A \rightarrow B \varepsilon^+ \mid \varepsilon^+ B \mid B) \in P} \varepsilon^+\text{-INT} \end{aligned}$$

To generate the sketch templates, we substitute two holes at an index-to-be-repaired, $H(\sigma, i) = \sigma_{1\dots i-1} \sigma_{i+1\dots n}$, then invoke the SAT solver. Five outcomes are then possible:

$$\sigma_1 \dots \sigma_{i-1} \text{ } \boxed{\gamma_1 \gamma_2} \text{ } \sigma_{i+1} \dots \sigma_n, \gamma_{1,2} = \varepsilon \quad (3)$$

$$\sigma_1 \dots \sigma_{i-1} \text{ } \boxed{\gamma_1 \gamma_2} \text{ } \sigma_{i+1} \dots \sigma_n, \gamma_1 \neq \sigma_i, \gamma_2 = \varepsilon \quad (4)$$

$$\sigma_1 \dots \sigma_{i-1} \text{ } \boxed{\gamma_1 \gamma_2} \text{ } \sigma_{i+1} \dots \sigma_n, \gamma_1 = \varepsilon, \gamma_2 \neq \sigma_i \quad (5)$$

$$\sigma_1 \dots \sigma_{i-1} \text{ } \boxed{\gamma_1 \gamma_2} \text{ } \sigma_{i+1} \dots \sigma_n, \gamma_1 = \sigma_i, \gamma_2 \neq \varepsilon \quad (6)$$

$$\sigma_1 \dots \sigma_{i-1} \text{ } \boxed{\gamma_1 \gamma_2} \text{ } \sigma_{i+1} \dots \sigma_n, \gamma_1 \notin \{\varepsilon, \sigma_i\}, \gamma_2 = \sigma_i \quad (7)$$

Eq. (3) corresponds to deletion, Eqs. (4, 5) correspond to substitution, and Eqs. (6, 7) correspond to insertion. The solutions returned by the SAT solver will be strictly equivalent to handling each edit operation as separate cases.

References

- [1] Valentin Antimirov. 1996. Partial derivatives of regular expressions and finite automaton constructions. *Theoretical Computer Science* 155, 2 (1996), 291–319.
- [2] Janusz A Brzozowski. 1964. Derivatives of regular expressions. *Journal of the ACM (JACM)* 11, 4 (1964), 481–494. http://maverick.uwaterloo.ca/reports/1964_JACM_Brzozowski.pdf
- [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>