

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

Anonymous Author(s)

## Abstract

Brzowski (1964) defines the derivative of a regular language as the suffixes that complete a known prefix. In this work, we establish a Galois connection with Valiant’s (1975) 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 ( $\Sigma$ ), 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, the CFG,  $P := \{S \rightarrow SS \mid (S) \mid ()\}$ , corresponds to the CNF:

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

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:

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

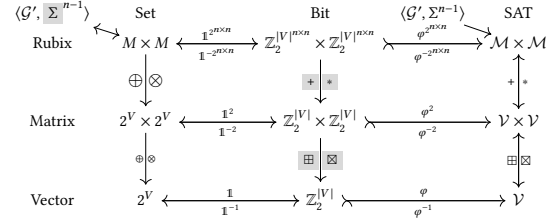
If we define  $\sigma_r^\dagger := \{w \mid (w \rightarrow \sigma_r) \in P\}$ , then initialize  $M_{r+1=c}^0(\mathcal{G}', e) := \sigma_r^\dagger$  and solve for the fixpoint  $M^* = M + M^2$ ,

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

we obtain the recognizer,  $R(\mathcal{G}', \sigma) := S \in \Lambda_\sigma^*? \Leftrightarrow \sigma \in \mathcal{L}(\mathcal{G})?$  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  $\omega$  is the least matrix multiplication upper bound (currently,  $\omega < 2.77$ ).

## 2 Method

Note that  $\bigoplus_{c=1}^n M_{r,c} \otimes M_{c,r}$  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 redefined as  $\boxplus, \boxtimes$  over bitvectors so the following diagram commutes,<sup>1</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,  $\phi^{-1}$  may not exist, as an arbitrary  $\mathcal{V}$  might take on zero, one, or in general, multiple values in  $\mathbb{Z}_2^{|V|}$ . Although holes may occur anywhere, let us consider two cases in which  $\Sigma^+$  is strictly left- or right-constrained, i.e.,  $xz, xz: \Sigma^{|x|+|z|}$ .

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.

Left Quotient

Right Quotient

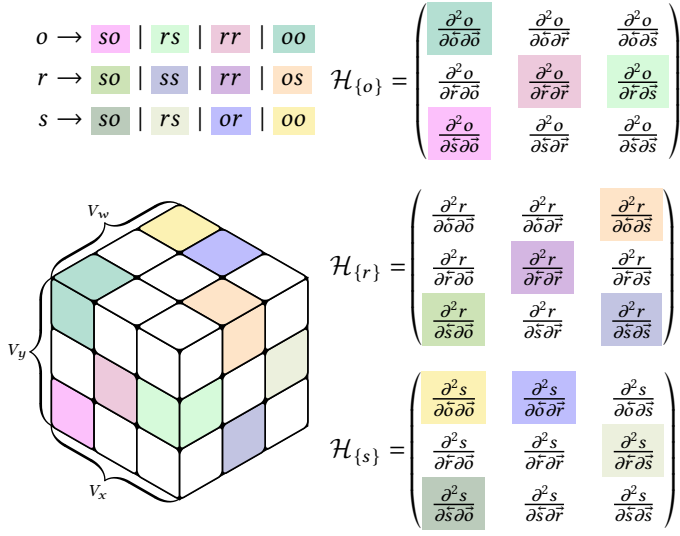
$$\frac{\partial}{\partial x} = \{ z \mid (w \rightarrow xz) \in P \}$$

$$\frac{\partial}{\partial z} = \{ x \mid (w \rightarrow xz) \in P \}$$



The left quotient coincides with the derivative operator first proposed by Brzowski [2] and Antimirov [1] over regular languages, lifted into the context-free setting (our work). When the root and LHS are fixed, e.g.,  $\frac{\partial S}{\partial x} : (\vec{V} \rightarrow S) \rightarrow \vec{V}$  returns the set of admissible nonterminals to the RHS. One may also consider a gradient operator,  $\vec{\nabla} S : (\vec{V} \rightarrow S) \rightarrow \vec{V}$ , which simultaneously tracks the partials with respect to a set of multiple LHS nonterminals produced by a fixed root.

<sup>1</sup>Hereinafter, we use gray highlighting to distinguish between expressions containing only **constants** from those which may contain free variables.



**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  $\mathcal{M}^*$ , we constrain the superposition of admissible parse forests.

## 2.1 Approximating the gradient

Given some unparseable string, i.e.,  $\sigma_1 \dots \sigma_n : \Sigma^n \cap \mathcal{L}(\mathcal{G})^c$ , where should we put holes to obtain  $\sigma' \in \mathcal{L}(\mathcal{G})$ ? We can estimate the effect of perturbing  $\sigma$  on  $\Lambda_\sigma^*$  by stochastically sampling *minibatches*  $\sigma \sim \Delta_k(\sigma)$  from the Levenshtein  $k$ -ball centered on  $\sigma$ , i.e., the space of all possible Levenshtein edits  $\leq k$ , and backpropagating  $\nabla S$  across upper-triangular entries of  $\mathcal{M}^*$ . In other words,

$$\sigma \sim \quad (2)$$

Let  $\mathbf{M} : \text{GF}(2^{n \times n})$  be a matrix whose structure is depicted in Eq. 3, where  $P$  is a feedback polynomial over  $\text{GF}(2^n)$  with coefficients  $P_{1..n}$  and semiring operators  $\oplus := \vee, \otimes := \wedge$ . Selecting any  $V \neq \mathbf{0}$  and coefficients  $P_{1..n}$  from a known *primitive polynomial* then powering the matrix  $\mathbf{M}$  generates an ergodic sequence over  $\text{GF}(2^n)$ , as shown in Eq. 4.

$$\mathbf{M}^t V = \begin{pmatrix} P_1 & \dots & P_n \\ \top & \circ & \circ \\ \circ & \circ & \circ \\ \circ & \circ & \top \end{pmatrix} \begin{pmatrix} V_1 \\ V_2 \\ V_3 \\ V_4 \end{pmatrix} \quad (3)$$

$$S = (V \quad MV \quad M^2V \quad M^3V \quad \dots \quad M^{2^n-1}V) \quad (4)$$

This sequence has *full periodicity*, in other words, for all  $i, j \in [0, 2^n)$ ,  $S_i = S_j \Rightarrow i = j$ . To uniformly sample  $\sigma \sim \Sigma^n$  without replacement, we form an injection  $\text{GF}(2^n) \rightarrow \Sigma^d$ ,

cycle through  $S$ , then discard samples that do not identify an element in any indexed dimension. This procedure rejects  $(1 - |\Sigma|2^{-\lceil \log_2 |\Sigma| \rceil})^d$  samples on average and requires  $\sim \mathcal{O}(1)$  per sample and  $\mathcal{O}(2^n)$  to exhaustively search the space.

## 2.2 Context-Sensitive Reachability

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

$$P_1 := \{ S \rightarrow LR, \quad L \rightarrow ab \mid aLb, \quad R \rightarrow c \mid cR \}$$

$$P_2 := \{ S \rightarrow LR, \quad R \rightarrow bc \mid bRc, \quad L \rightarrow a \mid aL \}$$

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. 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  $E(t_\cap, \mathcal{G}_i) = \{ w_i \mid (w_i \rightarrow t_\cap) \in P_i \}$  and glue them together:



**Figure 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 intersects orthonormal CNF unit productions  $w_i \rightarrow t_\cap$ , with outermost bitvector inhabitants representing  $S_i \in \Lambda_\sigma^*$ . Equations of this form are equiexpressive with the family of CSLs realizable by finite CFL intersection.

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

Although emptiness for CSLs is, in general, undecidable, unsatisfiability corresponds to emptiness of bounded CSLs. Anyhow, since we are working with bounded-width CSLs, everything collapses down to finite languages, which are always closed. Thus, Bar-Hillel and other elegant constructions become trivial when so restricted. So we can use it to decide impossible substrings and other decision problems which are typically intractable in more general settings.

## 3 Conclusion

This technique can be used to repair syntax errors in programming languages. It also has the nice advantage that it can provide partial parse trees for invalid strings, i.e., we can decode the tree from the structure of  $\mathbf{M}$  directly without any special error recovery.

## 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://maveric.uwaterloo.ca/reports/1964\\_JACM\\_Brzozowski.pdf](http://maveric.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>