# A Tree Sampler for Bounded Context-Free Languages

Breandan Considine

McGill University

breandan.considine@mail.mcgill.ca

## Abstract

The class of bounded context-free languages (BCFLs) consists of the subset of context-free languages which are finite. We provide a novel algorithm for sampling trees in BCFLs with and without replacement. Once the data structure is constructed, sampling trees is a straightforward matter of sampling integers uniformly without a replacement from a finite range. We demonstrate the utility of this technique on a dataset of Python statements.

## 1 Introduction

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

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

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

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

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

$$M^0 := \begin{pmatrix} \varnothing & \sigma_1^{\rightarrow} & \varnothing & \cdots & \varnothing \\ & & \ddots & & \vdots \\ & & & & \varnothing \\ & & & & \sigma_n^{\uparrow} \\ \varnothing & \cdots & \cdots & & \varnothing \end{pmatrix} \Rightarrow M^* = \begin{pmatrix} \varnothing & \sigma_1^{\rightarrow} & \Lambda & \cdots & \Lambda_\sigma^* \\ & & \ddots & & \Lambda \\ & & & & \sigma_n^{\uparrow} \\ \varnothing & \cdots & \cdots & & \varnothing \end{pmatrix}$$

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 [1], 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

We define the porous completion problem as follows:

**Definition 2.1** (Completion). Let $\underline{\Sigma} := \Sigma \cup \{\_\}$, where $\_$ represents a hole. We denote $\sqsubseteq \colon \Sigma^n \times \underline{\Sigma}^n$ as the relation $\{\langle \sigma', \sigma \rangle \mid \sigma_i \colon \Sigma \implies \sigma_i' = \sigma_i\}$ and the set $\{\sigma' \colon \Sigma^+ \mid \sigma' \sqsubseteq \sigma\}$ as $H(\sigma)$. Given $\sigma \colon \underline{\Sigma}^+$ we want to sample $\sigma' \sim H(\sigma) \cap \ell$.

$H(\sigma) \cap \ell$ is often a large-cardinality set, so we want a procedure which samples uniformly without replacement from the set, without enumerating the whole set and shuffling it.

We define an algebraic data type $\mathbb{T}_3 = (V \cup \Sigma) \rightharpoonup \mathbb{T}_2$ where $\mathbb{T}_2 = (V \cup \Sigma) \times (\mathbb{N} \rightharpoonup \mathbb{T}_2 \times \mathbb{T}_2)$[1]. Morally, we can think of $\mathbb{T}_2$ as an implicit set of possible trees sharing the same root, and $\mathbb{T}_3$ as a dictionary of possible $\mathbb{T}_2$ values indexed by possible roots, given by a specific CFG under a finite-length porous string. We construct $\hat{\sigma}_r = \dot{p}(\sigma_r)$ as follows:

$$\dot{p}(s \colon \Sigma) \mapsto \left\{ \mathbb{T}_2\big(w, \big[\langle \mathbb{T}_2(s), \mathbb{T}_2(\varepsilon) \rangle\big]\big) \mid (w \rightarrow s) \in P \right\}$$

$$\dot{p}(\_) \mapsto \bigoplus_{s \in \Sigma} p(s)$$

We then compute the fixpoint $M_\infty$ by redefining $\oplus, \otimes \colon \mathbb{T}_3 \times \mathbb{T}_3 \rightarrow \mathbb{T}_3$ as follows:

$$X \oplus Z \mapsto \bigcup_{k \in \pi_1(X \cup Z)} \left\{ k \Rightarrow \mathbb{T}_2(k, Q_x \cup Q_z) \mid Q_x \in \pi_2(X \circ k), Q_z \in \pi_2(Z \circ k) \right\}$$

$$X \otimes Z \mapsto \bigoplus_{w, x, z} \left\{ \mathbb{T}_2\big(w, \big[\langle X \circ x, Z \circ z \rangle\big]\big) \mid (w \rightarrow xz) \in P, x \in \pi_1(X), z \in \pi_1(Z) \right\}$$

Decoding trees from $(\Lambda_\sigma^* \circ S) \colon \mathbb{T}_2$ becomes a straightforward matter of enumeration using a recursive choice function that emits a sequence of binary trees generated by the CFG. We define this construction more precisely in § 2.1.

In our experiments, we provide a comparison of the performance of the SAT algebra and these two semirings, evaluated on a dataset of Python statements.
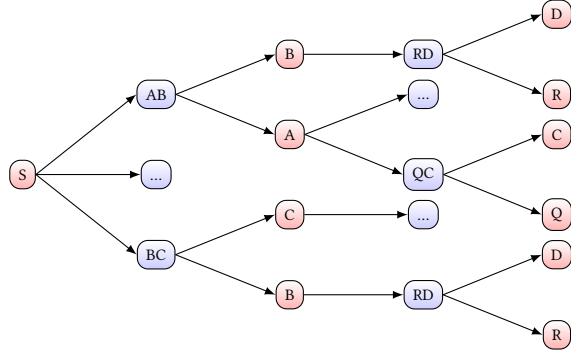
---

[1] Hereinafter, given a concrete $T \colon \mathbb{T}_2$, we shall refer to $\pi_1(T), \pi_2(T)$ as `root`$(T)$ and `children`$(T)$ respectively.

## 2.1 Pairing Breadth-Bounded Binary Trees

The type $\mathbb{T}_2$ of all possible trees that can be generated by a CFG in Chomksy Normal Form is identified by a recurrence relation:

$$L(p) = 1 + pL(p) \qquad P(a) = 1 + aL(P(a))^2$$

We depict it as a tree, where red nodes are `roots` and blue nodes are `children`:



The number of binary trees inhabiting a single instance of $\mathbb{T}_2$ is sensitive to the number of nonterminals and rule expansions in the grammar. To obtain the total number of trees with breadth $n$, we can take the intersection between a CFG and the regular language, $\mathcal{L}(G^\cap) := \mathcal{L}(\mathcal{G}) \cap \Sigma^n$, abstractly parse the string containing all holes, let $T = \Lambda^*_{\underline{\sigma}} \circ S$, and compute the total number of trees using the following recurrence:

$$|T : \mathbb{T}_2| \mapsto \begin{cases} 1 & \text{if } T \text{ is a leaf,} \\ \sum_{\langle T_1, T_2 \rangle \in \texttt{children}(T)} |T_1| \cdot |T_2| & \text{otherwise.} \end{cases}$$

To sample all trees in a given $T : \mathbb{T}_2$ uniformly without replacement, we first define a pairing function $\varphi : \mathbb{T}_2 \to \mathbb{Z}_{|T|} \to \texttt{BTree}$ as follows:

$$\varphi(T : \mathbb{T}_2, i : \mathbb{Z}_{|T|}) \mapsto \begin{cases} \left\langle \texttt{BTree}(\texttt{root}(T)), i \right\rangle & \text{if } T \text{ is a leaf,} \\[1em] \text{Let } b = |\texttt{children}(T)|, \\ \quad q_1, r = \left\langle \lfloor \frac{i}{b} \rfloor, i \pmod{b} \right\rangle, \\ \quad lb, rb = \texttt{children}[r], \\ \quad T_1, q_2 = \varphi(lb, q_1), \\ \quad T_2, q_3 = \varphi(rb, q_2) \text{ in} \\ \left\langle \texttt{BTree}(\texttt{root}(T), T_1, T_2), q_3 \right\rangle & \text{otherwise.} \end{cases}$$

Then, instead of sampling trees, we can simply sample integers uniformly without replacement from $\mathbb{Z}_{|T|}$ using the construction defined in § 2, and lazily decode them into trees.

## 3 Prior Work

Piantodosi define a similar construction, but it assumes the CFL is infinite and makes some additional assumptions about the CFG [3]. We provide a more general construction which works for any CFG. Sampling parse trees in CFGs can be viewed as sampling proofs in a limited kind of proof system [2].

## 4 Conclusion

We have presented a novel algorithm for sampling trees in bounded context-free languages with and without replacement. This technique has applications to code completion and program repair.

## References

[1] 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

[2] Andreas Opedal, Ran Zmigrod, Tim Vieira, Ryan Cotterell, and Jason Eisner. 2023. Efficient semiringweighted Earley parsing. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (ACL), Toronto, Canada.

[3] Steven T. Piantadosi. 2023. How to enumerate trees from a context-free grammar. arXiv:2305.00522 [cs.CL]

[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