

# A Tree Sampler for Bounded Context-Free Languages

Breandan Considine

## Main Idea

- If you can count it, you can sample it! Analytic combinatorics
- We implement a bijection between binary trees in bounded CFLs
- Allows for parallelizable replacement-free sampling from BCFLs

## Algebraic Parsing

Given a CFG  $\mathcal{G} := \langle V, \Sigma, P, S \rangle$  in Chomsky Normal Form (CNF), we may construct a recognizer  $R_{\mathcal{G}} : \Sigma^n \rightarrow \mathbb{B}$  for strings  $\sigma : \Sigma^n$  as follows. Let  $\in^V$  be our domain, where  $0$  is  $\emptyset$ ,  $\oplus$  is  $\cup$ , and  $\otimes$  be defined as:

$$s_1 \otimes s_2 := \{C \mid \langle A, B \rangle \in s_1 \times s_2, (C \rightarrow AB) \in P\}$$

Initializing  $\mathbf{M}_0[i, j](\mathcal{G}, \sigma) := \{A \mid i+1 = j, (A \rightarrow \sigma_i) \in P\}$  and searching for the least solution to  $\mathbf{M} = \mathbf{M} + \mathbf{M}^2$ , will produce a fixedpoint  $\mathbf{M}^*$ :

$$M^0 := \begin{pmatrix} \emptyset & \hat{\sigma}_1 & \emptyset & \emptyset \\ & & \emptyset & \\ & & & \hat{\sigma}_n \\ \emptyset & & & \emptyset \end{pmatrix} \Rightarrow M_{\infty} = \begin{pmatrix} \emptyset & \hat{\sigma}_1 & \Lambda & \Lambda_{\sigma}^* \\ & & \Lambda & \\ & & & \hat{\sigma}_n \\ \emptyset & & & \emptyset \end{pmatrix}$$

Valiant (1975) shows that  $\sigma \in \mathcal{L}(\mathcal{G})$  iff  $S \in \mathcal{T}$ , i.e.,  $1_{\mathcal{T}}(S) \iff 1_{\mathcal{L}(\mathcal{G})}(\sigma)$ .

## Parsing Dynamics

- The matrix  $\mathbf{M}_0$  is strictly upper triangular, i.e., nilpotent of degree  $n$
- The recognizer can be translated into a parser by storing *backpointers*

$\mathbf{M}_1 = \mathbf{M}_0 + \mathbf{M}_0^2$	$\mathbf{M}_2 = \mathbf{M}_1 + \mathbf{M}_1^2$	$\mathbf{M}_3 = \mathbf{M}_2 + \mathbf{M}_2^2 = \mathbf{M}_4$

- If we had a way to solve for  $\mathbf{M} = \mathbf{M} + \mathbf{M}^2$  directly, power iteration would be unnecessary and we could solve for  $\mathbf{M} = \mathbf{M}^2$  above the superdiagonal...

## Binarized CFL Sketching

- CYK parser can be lowered onto a Boolean tensor  $\mathbb{B}^{n \times n \times |V|}$  (Valiant, 1975)
- Binarized CYK parser can be compiled to SAT to solve for  $\mathbf{M}^*$  directly
- Enables sketch-based synthesis in either  $\sigma$  or  $\mathcal{G}$ : just use variables for holes!
- We simply encode the characteristic function, i.e.  $1_{\subseteq V} : V \rightarrow \mathbb{B}^{|V|}$
- $\oplus, \otimes$  are defined as  $\boxplus, \boxtimes$ , so that the following diagram commutes:

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

- These operators can be lifted into matrices and tensors in the usual way

## Method

We define an algebraic data type  $\mathbb{T}_3 = (V \cup \Sigma) \multimap \mathbb{T}_2$  where  $\mathbb{T}_2 = (V \cup \Sigma) \times (\mathbb{N} \multimap \mathbb{T}_2 \times \mathbb{T}_2)^a$ . 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  $\hat{\sigma}_r = \Lambda(\sigma_r)$  as follows:

$$\Lambda(s : \underline{\Sigma}) \mapsto \begin{cases} \bigoplus_{s \in \Sigma} \Lambda(s) & \text{if } s \text{ is a hole,} \\ \{\mathbb{T}_2(w, [\langle \mathbb{T}_2(s), \mathbb{T}_2(\varepsilon) \rangle]) \mid (w \rightarrow s) \in P\} & \text{otherwise.} \end{cases}$$

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

$$\begin{aligned} X \oplus Z &\mapsto \bigcup_{k \in \pi_1(X \cup Z)} \{k \Rightarrow \mathbb{T}_2(k, x \cup z) \mid x \in \pi_2(X \circ k), z \in \pi_2(Z \circ k)\} \\ X \otimes Z &\mapsto \bigoplus_{(w \rightarrow xz) \in P} \{\mathbb{T}_2(w, [\langle X \circ x, Z \circ z \rangle]) \mid x \in \pi_1(X), z \in \pi_1(Z)\} \end{aligned}$$

These operators group subtrees by their root nonterminal, then aggregate their children. Each  $\Lambda$  now becomes a dictionary indexed by the root nonterminal, which can be sampled by obtaining  $(\Lambda_{\sigma}^* \circ S) : \mathbb{T}_2$ , then recursively choosing *twins*.

<sup>a</sup>Given a  $T : \mathbb{T}_2$ , we may also refer to  $\pi_1(T), \pi_2(T)$  as  $\text{root}(T)$  and  $\text{children}(T)$  respectively, where children are pairs of conjoined twins.

## A Pairing Function for B-Trees

Let  $\mathbf{M} : \text{GF}(2^{n \times n})$  be a square matrix  $\mathbf{M}_{r,c}^0 = P_c$  if  $r = 0$  else  $1[c = r - 1]$ , where  $P$  is a feedback polynomial with coefficients  $P_{1..n}$  and  $\oplus := \vee, \otimes := \wedge$ :

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

Selecting any  $V \neq 0$  and coefficients  $P_j$  from a known *primitive polynomial*, then powering the matrix  $\mathbf{M}$  generates an ergodic sequence over  $\text{GF}(2^n)$ :

## Tidyparse IDE Plugin

```

let rec a = - - - - |
let rec filter p l = if p x then x :: - - - else - - -
let curry f = (fun x - - - - -)

S -> X
X -> A | V | ( X , X ) | X X | ( X )
A -> FUN | F | LI | M | L
FUN -> fun V -> X
F -> if X then X else X
M -> match V with Branch
Branch -> `|` X -> X | Branch Branch
L -> let V = X
L -> let rec V = X
LI -> L in X

V -> Vexp | ( Vexp ) | List | Vexp Vexp
Vexp -> Vname | FunName | Vexp V0 Vexp | B
Vexp -> ( Vname , Vname ) | Vexp Vexp | I
List -> [] | V :: V
Vname -> a | b | c | d | e | f | g | h | i
Vname -> j | k | l | m | n | o | p | q | r
Vname -> s | t | u | v | w | x | y | z
FunName -> foldright | map | filter
FunName -> curry | uncurry | ( V0 )
V0 -> + | - | * | / | >
V0 -> | < | `| | ' | &&
I -> 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
B -> true | false
    
```

