

Multi-dimensional Trees and a Chomsky-Schützenberger-Weir Representation Theorem for Simple Context-Free Tree Grammars

Makoto Kanazawa*

National Institute of Informatics, Tokyo, Japan
kanazawa@nii.ac.jp

Abstract. Weir [43] proved a Chomsky-Schützenberger-like representation theorem for the string languages of tree-adjoining grammars, where the Dyck language D_n in the Chomsky-Schützenberger characterization is replaced by the intersection $D_{2n} \cap g^{-1}(D_{2n})$, where g is a certain bijection on the alphabet consisting of $2n$ pairs of brackets. This paper presents a generalization of this theorem to the string languages generated by simple (i.e., linear and non-deleting) context-free tree grammars. This result is obtained through a natural generalization of the original Chomsky-Schützenberger theorem to the tree languages of simple context-free tree grammars. I use Baldwin and Strawn's [2] notion of multi-dimensional trees to state this latter theorem in a very general, abstract form.

Keywords: Context-free tree grammar; Multi-dimensional tree; Dyck language; Chomsky-Schützenberger theorem

1 Introduction

Weir [43] showed that every string language L generated by a tree-adjoining grammar [13] can be written as

$$L = h(R \cap D_{2n} \cap g^{-1}(D_{2n})),$$

where h is a homomorphism, R is a regular set, n is a positive integer, D_{2n} is the Dyck language over the alphabet Γ_{2n} consisting of $2n$ pairs of brackets $[_1,]_1, \dots, [_{2n},]_{2n}$, and g is the bijection on Γ_{2n} defined by

$$g([_{2i+1}) = [_{2i+1}, \quad g(]_{2i+1}) =]_{2i+2}, \quad g([_{2i+2}) =]_{2i+1}, \quad g(]_{2i+2}) = [_{2i+2},$$

* This work was in part supported by the Japan Society for the Promotion of Science Grant-in-Aid for Scientific Research (KAKENHI) (25330020). An earlier version of this paper is available as a technical report [17], where some of the missing details may be found. I am indebted to one of the anonymous reviewers for bringing Baldwin and Strawn's work to my attention.

for $i = 0, \dots, n-1$. The effect of the intersection with $g^{-1}(D_{2n})$ on the Dyck language D_{2n} is to make the consecutive odd-numbered and even-numbered brackets $[_{2i+1},]_{2i+1}, [_{2i+2},]_{2i+2}$ always appear as a group, in the configuration $[_{2i+1} [_{2i+2}]_{2i+2}]_{2i+1}$. When two such groups, say, $[_1,]_1, [_2,]_2$ and $[_3,]_3, [_4,]_4$, overlap, the only possible configurations are

$$\begin{aligned} &[_1 [_3 [_4]_4]_3]_2]_1, \\ &[_1 [_2]_2 [_3 [_4]_4]_3]_1, \\ &[_1 [_3 [_4 [_2]_2]_4]_3]_1, \\ &[_1 [_2 [_3 [_4]_4]_3]_2]_1, \end{aligned}$$

and those with the positions of the two groups interchanged. As Weir [43] showed, $D_{2n} \cap g^{-1}(D_{2n})$ is a non-context-free tree-adjoining language for every $n \geq 1$.

In this paper, I prove a generalization of Weir's theorem for *simple* (i.e., *linear* and *non-deleting*) *context-free tree grammars*¹ [35,9,22]: if L is the string language generated by a simple context-free tree grammar of rank $q-1$, then L can be written as

$$L = h(R \cap D_{qn} \cap g^{-1}(D_{qn})),$$

where h , R , and n are as before, D_{qn} is the Dyck language over the alphabet Γ_{qn} (containing qn pairs of brackets), and g is the bijection on Γ_{qn} defined by

$$\begin{aligned} g([_{qi+1}) &= [_{qi+1}, & g(]_{qi+1}) &=]_{qi+q}, \\ g([_{qi+j}) &=]_{qi+j-1}, & g(]_{qi+j}) &= [_{qi+j}, \end{aligned}$$

for $i = 0, \dots, n-1$ and $j = 2, \dots, q$. As with Weir's theorem, the intersection $D_{qn} \cap g^{-1}(D_{qn})$ is the string language of some simple context-free tree grammar of rank $q-1$. This result generalizes Weir's [43] because tree-adjoining grammars generate the same string languages as simple context-free tree grammars that are *monadic* (i.e., of rank 1) [30,12,23]. As in the original Chomsky-Schützenberger theorem [5,4], we can take R to be a *local* set and h to be *alphabetic* in the sense that h maps each symbol either to a symbol or to the empty string.

It is known [14,15] that the string languages of simple context-free tree grammars are exactly those generated by *multiple context-free grammars* [38] that are *well-nested* in the sense of [15]. For multiple context-free grammars in general, Yoshinaka et al. [45] have proved a Chomsky-Schützenberger-like representation theorem, but the analogy to the Chomsky-Schützenberger theorem is somewhat weak because their notion of a *multiple Dyck language* is given only by reference to a certain multiple context-free grammar, and does not seem to have other independent characterizations, analogous to the characterization of ordinary Dyck languages in terms of the cancellation law $[_i]_i \rightsquigarrow \varepsilon$. My result is obtained via a natural generalization of the Chomsky-Schützenberger theorem to the tree languages of simple context-free tree grammars, which may be of independent interest. This intermediate result is stated in terms of *Dyck tree languages*, which are exactly analogous to the original Dyck languages in that

¹ The term "simple context-free tree grammar" is taken from [8].

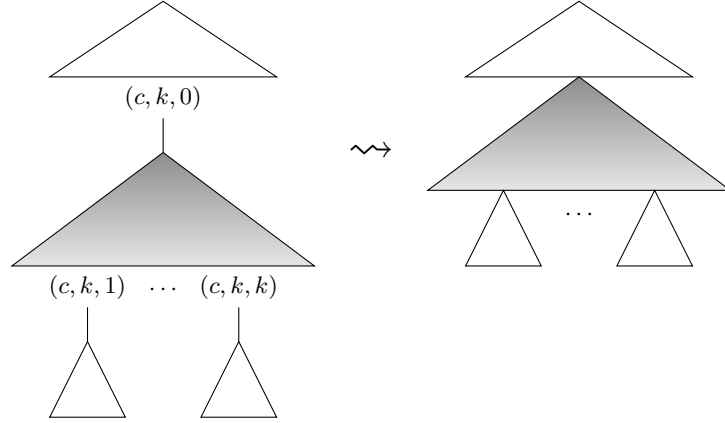


Fig. 1. A cancellation law for a Dyck tree language.

they have two equivalent definitions, one in terms of inductive definitions and one in terms of rewriting with cancellation laws.

Elements of a Dyck tree language are trees over an alphabet of the form $\tilde{\Sigma} = \Sigma \cup \{ (c, k, i) \mid c \in \Sigma, 0 \leq i \leq k \}$. For each c and k ,

$$(c, k, 0), (c, k, 1), \dots, (c, k, k)$$

form a matching group of symbols. A cancellation law for a Dyck tree language may be schematically depicted as in Fig. 1. The symbols $(c, k, 0), (c, k, 1), \dots, (c, k, k)$ may cancel out each other provided that the shaded region of the tree contains only symbols from Σ . A tree over $\tilde{\Sigma}$ belongs to the Dyck tree language if and only if successive applications of the cancellation law lead to a tree over Σ . I show that for every simple context-free tree language L , there exist an alphabet \mathcal{Y} , a local set R of trees over $\tilde{\mathcal{Y}}$, and a particularly simple kind of “linear and complete” tree homomorphism² h such that

$$L = h(R \cap DT_{\mathcal{Y}}),$$

where $DT_{\mathcal{Y}}$ is the Dyck tree language consisting of trees over $\tilde{\mathcal{Y}}$.

The intersection $D_{qn} \cap g^{-1}(D_{qn})$ in my generalization of Weir’s theorem comes from the “labeled bracketing” encoding of elements of a Dyck tree language, where the $k + 1$ pairs of brackets $\llbracket_{(c,k,0)}, \rrbracket_{(c,k,0)}, \llbracket_{(c,k,1)}, \rrbracket_{(c,k,1)}, \dots, \llbracket_{(c,k,k)}, \rrbracket_{(c,k,k)}$ form a group and always appear in the configuration

$$\llbracket_{(c,k,0)} \llbracket_{(c,k,1)} \llbracket_{(c,k,1)} \dots \llbracket_{(c,k,k)} \llbracket_{(c,k,k)} \rrbracket_{(c,k,0)}.$$

Simple context-free tree grammars hold interest for theoretical computational linguists because of their many attractive formal properties [15,14,19] and their

² More precisely, h is the composition of a projection and the operation of deleting certain unary-branching nodes.

ability to lexicalize tree-adjoining grammars without changing the set of derived trees [27]. The notion of a Dyck tree language I introduce in this paper also leads to a generalization of linear indexed grammars that is equivalent to simple context-free tree grammars in the same way that linear indexed grammars are equivalent to tree-adjoining grammars [18]. I believe these connections make it particularly interesting that the classic result of Chomsky and Schützenberger naturally extends to the level of simple context-free tree languages.

In order to emphasize the analogy between the string case and the tree case, I use the notion of a *multi-dimensional tree* introduced by Baldwin and Strawn [2] (and subsequently studied by Rogers [33,32] in connection with tree-adjoining grammars), and state many lemmas as general facts about m -dimensional trees. I use 3-dimensional trees to represent derivation trees of simple context-free tree grammars. Just as ordinary trees are encoded by strings of brackets belonging to an ordinary Dyck language, 3-dimensional trees are encoded by elements of a Dyck tree language. This correspondence extends to higher dimensions.³

2 Preliminaries

2.1 First-Child-Next-Sibling Encoding of Ordered Unranked Trees

In an ordered unranked tree, a node may have any number of children, and the children of the same node are linearly ordered. We do not consider unordered trees in this paper, so we call ordered unranked trees simply unranked trees. In the usual term notation for unranked trees [41], the unranked trees over a set Σ of labels are defined inductively as follows:

- If $c \in \Sigma$, then c is an unranked tree over Σ .
- If t_1, \dots, t_n are unranked trees over Σ ($n \geq 1$) and $c \in \Sigma$, then $c(t_1 \dots t_n)$ is an unranked tree over Σ .

There is a well-known way of encoding unranked trees into binary trees [24], often called the *first-child-next-sibling* encoding. We refer to a node in a binary tree by a string over the set $\{1, 2\}$. Thus, the set of nodes of a binary tree forms a prefix-closed subset T of $\{1, 2\}^*$. (Note that we do not assume binary trees to be *full* in the sense that each node has 0 or 2 children.) We write $u \cdot v$ for the concatenation of two strings $u, v \in \{1, 2\}^*$. In the first-child-next-sibling encoding of unranked trees, the relation $u \cdot 2 = v$ represents the relation “ v is the

³ When the present work was nearing completion, I learned of a recent paper by Sorokin [39], in which he states (without proof) a result similar to Theorem 40 below (Theorem 3 of [39]). (The statement of his theorem is actually closer to Lemma 36 below.) As will be clear to the reader, the emphasis of the present paper is very different from Sorokin’s. The merit of the present work lies not so much in Theorem 40 itself as in the method of obtaining it through a natural generalization of the constructions that can be used to prove the original Chomsky-Schützenberger theorem. (Sorokin’s own emphasis is on the use of monoid automata to characterize the string languages of simple context-free tree grammars.)

first child of u ", and the relation $u \cdot 1 = v$ represents the relation " v is the next sibling of u ". The child relation is then represented by the first-child relation composed with the reflexive transitive closure of the next-sibling relation. In this way, any non-empty finite prefix-closed subset T of $\{1, 2\}^*$ such that $1 \notin T$ encodes the set of nodes of some unranked tree. In general, an arbitrary non-empty finite prefix-closed subset of $\{1, 2\}^*$ encodes the nodes of a *hedge*, a finite, non-empty sequence of unranked trees.⁴ In this encoding, ε (empty string) is the root of the first tree, 1 is the root of the second tree, $1 \cdot 1$ is the root of the third tree, and so on.

Trees and hedges we consider in this paper are all labeled. Labeled unranked trees and hedges over Σ are represented by pairs of the form $\mathbf{T} = (T, \ell)$, where T is a non-empty finite prefix-closed subset of $\{1, 2\}^*$ and ℓ is a function from T to Σ . The set of labeled unranked trees over Σ is denoted \mathbb{T}_Σ .

2.2 Dyck Languages

For $n \geq 1$, let $\Gamma_n = \bigcup_{i=1}^n \{[i,]_i\}$. For each i , the two symbols $[i,]_i$ are regarded as a matching pair of brackets. Define a binary relation \rightsquigarrow on Γ_n^* by

$$\rightsquigarrow = \{ (u [i]_i v, uv) \mid u, v \in \Gamma_n^*, 1 \leq i \leq n \}.$$

The *Dyck language* D_n is defined by

$$D_n = \{ v \in \Gamma_n^* \mid v \rightsquigarrow^* \varepsilon \},$$

where \rightsquigarrow^* denotes the reflexive transitive closure of the relation \rightsquigarrow . An alternative way of defining D_n is by the following context-free grammar:

$$\begin{aligned} S &\rightarrow \varepsilon \mid AS, \\ A &\rightarrow [{}_1 S]_1 \mid \cdots \mid [{}_n S]_n. \end{aligned}$$

The set D'_n of *Dyck primes* is defined by

$$D'_n = (D_n - \{\varepsilon\}) - (D_n - \{\varepsilon\})^2.$$

Alternatively, the set D'_n is defined by the nonterminal A in the above context-free grammar.

Unranked trees and hedges can be represented by elements of Dyck languages. If Σ is a set of symbols, let

$$\Gamma_\Sigma = \bigcup_{c \in \Sigma} \{[c,]_c\}.$$

⁴ Sometimes the empty sequence of unranked trees is also allowed as a hedge, but we exclude it here in order to be able to encode all hedges into binary trees. Note that Knuth [24], Takahashi [40], and Baldwin and Strawn [2] used *forest* instead of *hedge*, the term I adopt here following [6].

We write D_Σ and D'_Σ for the Dyck language and the set of Dyck primes over this alphabet. Using the standard term notation for labeled unranked trees, define the *string encoding* function **enc** from labeled unranked trees and hedges over Σ to strings over Γ_Σ by

$$\begin{aligned} \mathbf{enc}(c) &= \llbracket_c \rrbracket_c, \\ \mathbf{enc}(c(t_1 \dots t_n)) &= \llbracket_c \mathbf{enc}(t_1 \dots t_n) \rrbracket_c, \\ \mathbf{enc}(t_1 \dots t_n) &= \mathbf{enc}(t_1) \mathbf{enc}(t_2 \dots t_n) \quad \text{for } n \geq 2. \end{aligned}$$

It is clear that the function **enc** maps any unranked tree over Σ to an element of D'_Σ , and any hedge over Σ to an element of $D_\Sigma - \{\varepsilon\}$. Conversely, it is easy to see that any element of D'_Σ encodes a tree over Σ , and any element of $D_\Sigma - \{\varepsilon\}$ encodes a hedge over Σ . These correspondences are bijections.

2.3 Context-Free Tree Grammars

We deviate from the standard practice and let a context-free tree grammar generate a set of unranked trees. Thus, the terminal alphabet of a context-free tree grammar will be unranked. In contrast, the set of nonterminals will be a ranked alphabet, as in the standard definition.

A *ranked alphabet* is a union $\mathcal{Y} = \bigcup_{r \in \mathbb{N}} \mathcal{Y}^{(r)}$ of disjoint sets of symbols. If $f \in \mathcal{Y}^{(r)}$, r is the *rank* of f . If Σ is an (unranked) alphabet and \mathcal{Y} a ranked alphabet ($\Sigma \cap \mathcal{Y} = \emptyset$), let $\mathbb{T}_{\Sigma, \mathcal{Y}}$ be the set of trees $\mathbf{T} \in \mathbb{T}_{\Sigma \cup \mathcal{Y}}$ such that whenever a node of \mathbf{T} is labeled by some $f \in \mathcal{Y}$, then the number of its children is equal to the rank of f .

For convenience, we use the term representation of trees. The set $\mathbb{T}_{\Sigma, \mathcal{Y}}$ can be defined inductively as follows:

1. If $f \in \Sigma \cup \mathcal{Y}^{(0)}$, then $f \in \mathbb{T}_{\Sigma, \mathcal{Y}}$;
2. If $f \in \Sigma \cup \mathcal{Y}^{(n)}$ and $t_1, \dots, t_n \in \mathbb{T}_{\Sigma, \mathcal{Y}}$ ($n \geq 1$), then $f(t_1 \dots t_n) \in \mathbb{T}_{\Sigma, \mathcal{Y}}$.

In order to define the notion of a context-free tree grammar, we need a countably infinite supply of variables x_1, x_2, x_3, \dots . The set consisting of the first n variables is denoted X_n (i.e., $X_n = \{x_1, \dots, x_n\}$). The notation $\mathbb{T}_{\Sigma, \mathcal{Y}}(X_n)$ denotes the set $\mathbb{T}_{\Sigma, \mathcal{Y} \cup X_n}$, where members of X_n are all assumed to have rank 0. A tree in $\mathbb{T}_{\Sigma, \mathcal{Y}}(X_n)$ is often written $t[x_1, \dots, x_n]$, displaying the variables. If $t[x_1, \dots, x_n] \in \mathbb{T}_{\Sigma, \mathcal{Y}}(X_n)$ and $t_1, \dots, t_n \in \mathbb{T}_{\Sigma, \mathcal{Y}}$, then $t[t_1, \dots, t_n]$ denotes the result of substituting t_1, \dots, t_n for x_1, \dots, x_n , respectively, in $t[x_1, \dots, x_n]$. An element $t[x_1, \dots, x_n]$ of $\mathbb{T}_{\Sigma, \mathcal{Y}}(X_n)$ is an *n-context* if for each $i = 1, \dots, n$, x_i occurs exactly once in $t[x_1, \dots, x_n]$. (In the literature, an *n-context* is sometimes called a *simple tree*.)

A *context-free tree grammar* [35,9] is a quadruple $G = (N, \Sigma, P, S)$, where

1. N is a finite ranked alphabet of nonterminals,
2. Σ is a finite unranked alphabet of terminals,
3. S is a nonterminal of rank 0, and

4. P is a finite set of productions of the form

$$B(x_1 \dots x_n) \rightarrow t[x_1, \dots, x_n],$$

where $B \in N^{(n)}$ and $t[x_1, \dots, x_n] \in \mathbb{T}_{\Sigma, N}(X_n)$.

The *rank* of G is $\max\{r \mid N^{(r)} \neq \emptyset\}$.

For every $s, s' \in \mathbb{T}_{\Sigma, N}$, $s \Rightarrow_G s'$ is defined to hold if and only if there is a 1-context $c[x_1] \in \mathbb{T}_{\Sigma, N}(X_1)$, a production $B(x_1 \dots x_n) \rightarrow t[x_1, \dots, x_n]$ in P , and trees $t_1, \dots, t_n \in \mathbb{T}_{\Sigma, N}$ such that

$$\begin{aligned} s &= c[B(t_1 \dots t_n)], \\ s' &= c[t[t_1, \dots, t_n]]. \end{aligned}$$

The relation \Rightarrow_G^* on $\mathbb{T}_{\Sigma, N}$ is defined as the reflexive transitive closure of \Rightarrow_G . The *tree language* generated by a context-free tree grammar G , denoted by $L(G)$, is defined as follows:

$$L(G) = \{t \in \mathbb{T}_{\Sigma} \mid S \Rightarrow_G^* t\}.$$

The *string language* generated by G is

$$\mathbf{y}(L(G)) = \{\mathbf{y}(t) \mid t \in L(G)\},$$

where $\mathbf{y}(t)$ is the yield of t in the usual sense.

A context-free tree grammar $G = (N, \Sigma, P, S)$ is said to be *simple* if for every production

$$B(x_1 \dots x_n) \rightarrow t[x_1, \dots, x_n]$$

in P , $t[x_1, \dots, x_n]$ is an n -context. We let $\text{CFT}_{\text{sp}}(r)$ stand for the family of tree languages L such that $L = L(G)$ for some simple context-free tree grammar G whose rank does not exceed r . We write $\mathbf{y}\text{CFT}_{\text{sp}}(r)$ for the corresponding string languages $\{\mathbf{y}(L) \mid L \in \text{CFT}_{\text{sp}}(r)\}$.

Let $\{y_1, \dots, y_k\}$ be a ranked alphabet, where for $i = 1, \dots, k$, r_i is the rank of y_i . Let $t_i[x_1, \dots, x_{r_i}]$ be an r_i -context. For a tree $t \in \mathbb{T}_{\Sigma, \{y_1, \dots, y_k\}}(X_n)$, we define $t[t_i[x_1, \dots, x_{r_i}]/y_i]$ inductively as follows:

$$\begin{aligned} c(u_1 \dots u_m)[t_i[x_1, \dots, x_{r_i}]/y_i] &= c(u_1[t_i[x_1, \dots, x_{r_i}]/y_i] \dots u_m[t_i[x_1, \dots, x_{r_i}]/y_i]) \\ &\text{if } c \in \Sigma, \end{aligned}$$

$$x_j[t_i[x_1, \dots, x_{r_i}]/y_i] = x_j,$$

$$y_i(u_1 \dots u_{r_i})[t_i[x_1, \dots, x_{r_i}]/y_i] = t_i[u_1[t_i[x_1, \dots, x_{r_i}]/y_i], \dots, u_{r_i}[t_i[x_1, \dots, x_{r_i}]/y_i]],$$

$$\begin{aligned} y_j(u_1 \dots u_{r_j})[t_i[x_1, \dots, x_{r_i}]/y_i] &= y_j(u_1[t_i[x_1, \dots, x_{r_i}]/y_i] \dots u_{r_j}[t_i[x_1, \dots, x_{r_i}]/y_i]) \\ &\text{if } j \neq i. \end{aligned}$$

(Here, the notation $c(u_1 \dots u_m)$ stands for c when $m = 0$, and likewise with $y_j(u_1 \dots u_{r_j})$.)

Let $G = (N, \Sigma, P, S)$ be a simple context-free tree grammar. The *derivation trees* of G and their *tree yield* are defined inductively as follows:

- Let $\pi = B(x_1 \dots x_n) \rightarrow t[x_1, \dots, x_n]$ be a production in P with no nonterminal occurring in $t[x_1, \dots, x_n]$. Then $\mathbf{d} = \pi$ is a derivation tree of sort B and its tree yield is $\mathbf{ty}(\mathbf{d}) = t[x_1, \dots, x_n]$.
- Let $\pi = B(x_1 \dots x_n) \rightarrow t[x_1, \dots, x_n]$ be a production in P with at least one nonterminal occurring in $t[x_1, \dots, x_n]$. Let v_1, \dots, v_k be the pre-order listing of the nodes of $t[x_1, \dots, x_n]$ labeled by nonterminals. Let B_i be the label of v_i , for $i = 1, \dots, k$. If \mathbf{d}_i is a derivation tree of sort B_i for $i = 1, \dots, k$, then

$$\pi(\mathbf{d}_1 \dots \mathbf{d}_k)$$

is a derivation tree of sort B . If $\bar{t}[x_1, \dots, x_n]$ is a tree just like $t[x_1, \dots, x_n]$ except that the label of v_i is changed to y_i , where y_1, \dots, y_k are new symbols, then the tree yield $\mathbf{ty}(\mathbf{d})$ of $\mathbf{d} = \pi(\mathbf{d}_1 \dots \mathbf{d}_k)$ is defined by

$$\mathbf{ty}(\mathbf{d}) = \bar{t}[x_1, \dots, x_n][\mathbf{ty}(\mathbf{d}_1)/y_1] \dots [\mathbf{ty}(\mathbf{d}_k)/y_k].$$

Note that if \mathbf{d} is a derivation tree of sort B and n is the rank of B , then $\mathbf{ty}(\mathbf{d})$ is an n -context, so the right-hand side of the above equation is well-defined if y_i is regarded as a symbol whose rank equals the rank of B_i . It is well known that if $G = (N, \Sigma, P, S)$ is a simple context-free grammar, then

$$L(G) = \{ \mathbf{ty}(\mathbf{d}) \mid \mathbf{d} \text{ is a derivation tree of } G \text{ of sort } S \}.$$

Example 1. Consider a simple context-free tree grammar $G = (N, \Sigma, P, S)$, where $N = N^{(0)} \cup N^{(2)} = \{S\} \cup \{B\}$, $\Sigma = \{a_1, a_2, a_3, a_4, a_5, a_6, e, g, h\}$, and P consists of the following rules:

$$\begin{aligned} \pi_1: S &\rightarrow B(ee), \\ \pi_2: B(x_1x_2) &\rightarrow h(a_1B(h(a_2x_1a_3)h(a_4x_2a_5))a_6), \\ \pi_3: B(x_1x_2) &\rightarrow g(x_1x_2). \end{aligned}$$

The following trees are derivation trees of this grammar:

$$\pi_1(\pi_2(\pi_3)) \quad \pi_1(\pi_2(\pi_2(\pi_3)))$$

We have

$$\begin{aligned} \mathbf{ty}(\pi_3) &= g(x_1x_2), \\ \mathbf{ty}(\pi_2(\pi_3)) &= h(a_1g(h(a_2x_1a_3)h(a_4x_2a_5))a_6), \\ \mathbf{ty}(\pi_1(\pi_2(\pi_3))) &= h(a_1g(h(a_2ea_3)h(a_4ea_5))a_6), \\ \mathbf{ty}(\pi_2(\pi_2(\pi_3))) &= h(a_1h(a_1g(h(a_2h(a_2x_1a_3)a_3)h(a_4h(a_4x_2a_5)a_5))a_6)a_6), \\ \mathbf{ty}(\pi_1(\pi_2(\pi_2(\pi_3)))) &= h(a_1h(a_1g(h(a_2h(a_2ea_3)a_3)h(a_4h(a_4ea_5)a_5))a_6)a_6). \end{aligned}$$

The string language generated by this grammar is

$$\mathbf{y}(L(G)) = \{ a_1^n a_2^n e a_3^n a_4^n e a_5^n a_6^n \mid n \geq 0 \}.$$

The string languages of simple context-free grammars are the languages generated by *non-duplicating macro grammars* [11], studied by Seki and Kato [37].⁵ They also coincide with the languages generated by *well-nested multiple context-free grammars* [15].

3 The Chomsky-Schützenberger Theorem

There are many different proofs of the Chomsky-Schützenberger theorem for context-free languages offered in the literature. Here, I give a proof based on the relation between context-free languages and local sets of unranked trees. In this proof, a subset of a Dyck language may be viewed as an alternative representation of a given local set of unranked trees (or the set of derivation trees of a given context-free grammar).⁶ This nicely captures Chomsky’s [5] original concern, and will serve as a starting point for our generalization of the theorem to simple context-free tree grammars. For simplicity and uniformity with the case of tree languages, we restrict ourselves to ε -free context-free languages here. It is of course easy to lift this restriction.

Let $\mathbf{T} = (T, \ell)$ be an unranked tree (in first-child-next-sibling encoding). We define three binary relations on T :⁷

$$\begin{aligned} \prec_2^T &= \{ (u, v) \in T \times T \mid u \cdot 2 = v \}, \\ \prec_1^T &= \{ (u, v) \in T \times T \mid u \cdot 1 = v \}, \\ \prec^T &= \{ (u, v) \in T \times T \mid u \prec_2^T \circ (\prec_1^T)^* v \}. \end{aligned}$$

The relation \prec^T is the child relation on the nodes of \mathbf{T} .

If $A, Z \subseteq \Sigma$ and $I \subseteq \Sigma \times \Sigma^+$ are finite sets, define $\text{Loc}(A, Z, I)$ to be the set of all trees $\mathbf{T} = (T, \ell)$ in \mathbb{T}_Σ that satisfy the following conditions:

- L1. $\ell(\varepsilon) \in A$.
- L2. $u \in T - \text{dom}(\prec_2^T)$ implies $\ell(u) \in Z$.
- L3. $u \prec_2^T v_1 \prec_1^T \dots \prec_1^T v_n \notin \text{dom}(\prec_1^T)$ ($n \geq 1$) implies $(\ell(u), \ell(v_1) \dots \ell(v_n)) \in I$.

A set $L \subseteq \mathbb{T}_\Sigma$ is *local* [42,40] if there are finite sets $A, Z \subseteq \Sigma$ and $I \subseteq \Sigma \times \Sigma^+$ such that $L = \text{Loc}(A, Z, I)$.

We introduce another notion of locality, which is in a sense more natural from the point of view of first-child-next-sibling encoding. If $A, Z, Y \subseteq \Sigma$ and $K, J \subseteq \Sigma \times \Sigma$, define $\text{SLoc}(A, Z, K, Y, J)$ to be the set of all trees $\mathbf{T} = (T, \ell)$ in \mathbb{T}_Σ that satisfy the following conditions:

⁵ At the level of string languages, simple context-free tree grammars correspond to non-duplicating and argument-preserving (i.e., non-deleting) macro grammars, which are equivalent to non-duplicating macro grammars (Lemma 7 of [37]). Seki and Kato [37] called non-duplicating macro grammars *variable-linear*.

⁶ Among the proofs found in well-known textbooks, the one closest to the present proof seems to be the one given by Kozen [25].

⁷ When R and S are binary relations on some set, we write $R \circ S$ for the composition of R and S , and write R^* for the reflexive transitive closure of R .

- SL1. $\ell(\varepsilon) \in A$.
 SL2. $u \in T - \text{dom}(\prec_2^T)$ implies $\ell(u) \in Z$.
 SL3. $u \prec_2^T v$ implies $(\ell(u), \ell(v)) \in K$.
 SL4. $u \neq \varepsilon$ and $u \in T - \text{dom}(\prec_1^T)$ imply $\ell(u) \in Y$.
 SL5. $u \prec_1^T v$ implies $(\ell(u), \ell(v)) \in J$.

We call $L \subseteq \mathbb{T}_\Sigma$ *super-local* if there are finite sets $A, Z, Y \subseteq \Sigma$ and $K, J \subseteq \Sigma \times \Sigma$ such that $L = \text{SLoc}(A, Z, K, Y, J)$.⁸

A set of strings $L \subseteq \Sigma^+$ is *local*⁹ if there are finite sets $A, Z \subseteq \Sigma$ and $I \in \Sigma^2$ such that

$$L = A\Sigma^* \cap \Sigma^*Z - (\Sigma^*(\Sigma^2 - I)\Sigma^*).$$

In this paper, we allow the alphabet Σ to be infinite, but any local subset of Σ^+ is included in Σ_0^+ for some finite subset Σ_0 of Σ ; likewise, any local or super-local subset of \mathbb{T}_Σ is included in \mathbb{T}_{Σ_0} for some finite $\Sigma_0 \subseteq \Sigma$.

We extend the string encoding function \mathbf{enc} to a function from $\mathcal{P}(\mathbb{T}_\Sigma)$ to $\mathcal{P}(\Gamma_\Sigma^+)$ by $\mathbf{enc}(L) = \{\mathbf{enc}(\mathbf{T}) \mid \mathbf{T} \in L\}$. In general, when $L \subseteq \mathbb{T}_\Sigma$ is local, there may be no regular set $L' \subseteq \Gamma_\Sigma^+$ such that $\mathbf{enc}(L) = L' \cap D'_\Sigma$ [26,16].

Lemma 2. *Let $L \subseteq \mathbb{T}_\Sigma$. If L is super-local, then there is a local set of strings $L' \subseteq \Gamma_\Sigma^+$ such that $\mathbf{enc}(L) = L' \cap D'_\Sigma$.*

Proof. Suppose that $A, Z, Y \subseteq \Sigma$ and $K, J \subseteq \Sigma \times \Sigma$ are finite sets such that $L = \text{SLoc}(A, Z, K, Y, J)$. Let

$$\begin{aligned} A' &= \{ \llbracket_c \mid c \in A \}, \\ Z' &= \{ \rrbracket_c \mid c \in A \}, \\ I &= \{ \llbracket_c \llbracket_d \mid (c, d) \in K \} \cup \{ \llbracket_c \rrbracket_c \mid c \in Z \} \cup \{ \rrbracket_c \llbracket_d \mid (c, d) \in J \} \cup \\ &\quad \{ \rrbracket_c \rrbracket_d \mid c \in Y, d \in \Sigma \}. \end{aligned}$$

Let $L' \subseteq \Gamma_\Sigma^+$ be the local set of strings defined by

$$L' = A' \Gamma_\Sigma^* \cap \Gamma_\Sigma^* Z' - (\Gamma_\Sigma^*(\Gamma_\Sigma^2 - I) \Gamma_\Sigma^*).$$

It is straightforward to show that $\mathbf{enc}(L) = L' \cap D'_\Sigma$. □

Note that the converse of the above lemma does not necessarily hold, because L' can place a restriction on \rrbracket_d that can follow \rrbracket_c . For example, $L = \{a(bc), a(bd), e(bc)\}$ is not super-local, even though $\mathbf{enc}(L)$ is local.

A mapping $\pi: \Sigma \rightarrow \Sigma'$ is called a *projection*. A projection π is extended to a function from \mathbb{T}_Σ to $\mathbb{T}_{\Sigma'}$ and to a function from $\mathcal{P}(\mathbb{T}_\Sigma)$ to $\mathcal{P}(\mathbb{T}_{\Sigma'})$ in obvious ways.

⁸ This notion of super-locality was called \tilde{F} -locality by Takahashi [40].

⁹ This notion of a local set is slightly different from McNaughton and Papert's [29] notion of a *strictly 2-testable* language. In the literature, a local set of strings is sometimes called *strictly 2-local* (for example, [34]). Eilenberg [7], Takahashi [40], and Perrin [31] use "local" in the present sense. Local sets of strings were called "standard regular events" by Chomsky and Schützenberger [4].

Lemma 3. *Let $L \subseteq \mathbb{T}_\Sigma$ be a local set. There exist a finite alphabet Σ' , a super-local set $L' \subseteq \mathbb{T}_{\Sigma'}$, and a projection $\pi: \Sigma' \rightarrow \Sigma$ such that $L = \pi(L')$. Moreover, π maps L' bijectively to L .*

Proof. Let $\mathbf{T} \in L$. We change the label of each node v of \mathbf{T} by a pair $(c_1 \dots c_n, i)$, where $c_1 \dots c_n$ is the string of labels c_1, \dots, c_n of the siblings of v , including v itself, in the order from left to right, and i is the position of v among its siblings. The relabeled trees obtained this way form a super-local set, and we can get back the original trees by a projection.

Formally, let¹⁰

$$\Sigma'' = \{ (w, i) \mid w \in \Sigma^+, 1 \leq i \leq |w| \},$$

and define a projection $\pi: \Sigma'' \rightarrow \Sigma$ by

$$\pi((c_1 \dots c_n), i) = c_i.$$

Suppose that $A, Z \subseteq \Sigma$ and $I \in \Sigma \times \Sigma^+$ are finite sets such that $L = \text{Loc}(A, Z, I)$. Let

$$\begin{aligned} F &= A \cup \{ w \in \Sigma^+ \mid (c, w) \in I \}, \\ \Sigma' &= \{ (w, i) \in \Sigma'' \mid w \in F \}. \end{aligned}$$

Note that Σ' is a finite subset of Σ'' . Let

$$\begin{aligned} A' &= \{ (c, 1) \mid c \in A \}, \\ Z' &= \{ (c_1 \dots c_n, i) \in \Sigma' \mid c_i \in Z \}, \\ K &= \{ ((d_1 \dots d_l, i), (c_1 \dots c_n, 1)) \mid (d_1 \dots d_l, i) \in \Sigma', (d_i, c_1 \dots c_n) \in I \} \\ Y &= \{ (c_1 \dots c_n, i) \in \Sigma' \mid i = n \}, \\ J &= \{ ((c_1 \dots c_n, i), (c_1 \dots c_n, i + 1)) \mid (c_1 \dots c_n, i) \in \Sigma', i \leq n - 1 \}. \end{aligned}$$

These sets are all finite. Let $L' \subseteq \mathbb{T}_{\Sigma''}$ be the super-local set defined by $L' = \text{SLoc}(A', Z', K, Y, J)$. Clearly, $L' \subseteq \mathbb{T}_{\Sigma'}$. We show that L' and π (restricted to Σ') satisfy the required properties.

For each $\mathbf{T} = (T, \ell^{\mathbf{T}}) \in \mathbb{T}_\Sigma$, define a tree $\hat{\mathbf{T}} = (T, \ell^{\hat{\mathbf{T}}}) \in \mathbb{T}_{\Sigma''}$ by

$$\ell^{\hat{\mathbf{T}}}(\varepsilon) = (\ell^{\mathbf{T}}(\varepsilon), 1), \tag{1}$$

$$\begin{aligned} \ell^{\hat{\mathbf{T}}}(u \cdot 2 \cdot 1^{i-1}) &= (\ell^{\mathbf{T}}(u \cdot 2) \ell^{\mathbf{T}}(u \cdot 2 \cdot 1) \dots \ell^{\mathbf{T}}(u \cdot 2 \cdot 1^{n-1}), i) \\ &\text{if } u \cdot 2 \cdot 1^{n-1} \in T - \text{dom}(\prec_1^T) \text{ and } 1 \leq i \leq n. \end{aligned} \tag{2}$$

It is clear that $\pi(\hat{\mathbf{T}}) = \mathbf{T}$ for all $\mathbf{T} \in \mathbb{T}_\Sigma$. Our goal is to show

$$L' = \{ \hat{\mathbf{T}} \mid \mathbf{T} \in L \}.$$

¹⁰ If w is a string, we write $|w|$ for the length of w . We use $|\cdot|$ both for the length of a string and for the cardinality of a set. The context should make it clear which is intended.

This clearly implies that π is a bijection from L' to L .

We begin by showing that for all $\mathbf{T} \in \mathbb{T}_\Sigma$,

$$\mathbf{T} \in L \text{ if and only if } \hat{\mathbf{T}} \in L'. \quad (3)$$

This follows from five observations. Firstly, note the following:

- Suppose $u \in T - \text{dom}(\prec_1^T)$. Then $\ell^{\hat{\mathbf{T}}}(u)$ is of the form $(c_1 \dots c_n, n)$, which means that $\ell^{\hat{\mathbf{T}}}(u) \in Y$ if $\ell^{\hat{\mathbf{T}}}(u) \in \Sigma'$.
- Suppose $u \prec_1^T v$. Then $(\ell^{\hat{\mathbf{T}}}(u), \ell^{\hat{\mathbf{T}}}(v))$ is of the form $((c_1 \dots c_n, i), (c_1 \dots c_n, i+1))$, which means that $(\ell^{\hat{\mathbf{T}}}(u), \ell^{\hat{\mathbf{T}}}(v)) \in J$ if $\ell^{\hat{\mathbf{T}}}(u) \in \Sigma'$.

Thus, $\hat{\mathbf{T}}$ satisfies the last two conditions SL4 and SL5 for membership in $\text{SLoc}(A', Z', K, Y, J)$ whenever $\hat{\mathbf{T}} \in \mathbb{T}_{\Sigma'}$. Secondly, the following biconditional always holds:

- $\ell^{\mathbf{T}}(\varepsilon) \in A$ if and only if $\ell^{\hat{\mathbf{T}}}(\varepsilon) \in A'$.

Thirdly, the following biconditional holds whenever $\ell^{\mathbf{T}}(v) \in \Sigma'$:

- $\ell^{\mathbf{T}}(v) \in Z$ if and only if $\ell^{\hat{\mathbf{T}}}(v) \in Z'$.

Fourthly, if $u \cdot 2 \cdot 1^{n-1} \in T - \text{dom}(\prec_1^T)$ and $\ell^{\hat{\mathbf{T}}}(u) \in \Sigma'$, then the following biconditional holds:

- $(\ell^{\mathbf{T}}(u), \ell^{\mathbf{T}}(u \cdot 2) \ell^{\mathbf{T}}(u \cdot 2 \cdot 1) \dots \ell^{\mathbf{T}}(u \cdot 2 \cdot 1^{n-1})) \in I$ if and only if $(\ell^{\hat{\mathbf{T}}}(u), \ell^{\hat{\mathbf{T}}}(u \cdot 2)) \in K$.

Lastly, it is easy to see that $\mathbf{T} \in L$ implies $\hat{\mathbf{T}} \in \mathbb{T}_{\Sigma'}$. Combining these five observations, we get (3).

It follows from the “only if” direction of (3) that $\{\hat{\mathbf{T}} \mid \mathbf{T} \in L\} \subseteq L'$. To establish the converse inclusion, we show that

$$\text{if } \mathbf{T}' \in L' \text{ and } \mathbf{T} = \pi(\mathbf{T}'), \text{ then } \mathbf{T}' = \hat{\mathbf{T}}.$$

This together with the “if” direction of (3) clearly implies $L' \subseteq \{\hat{\mathbf{T}} \mid \mathbf{T} \in L\}$.

So suppose $\mathbf{T}' = (T, \ell^{\mathbf{T}'}) \in L'$, and let $\mathbf{T} = (T, \ell^{\mathbf{T}}) = \pi(\mathbf{T}')$. All we need to show is that the equations (1) and (2) hold with \mathbf{T}' in place of $\hat{\mathbf{T}}$. As for (1), it follows from the fact that $\ell^{\mathbf{T}'}(\varepsilon) \in A'$. As for (2), suppose $u \cdot 2 \cdot 1^{n-1} \in T - \text{dom}(\prec_1^T)$. Since $(\ell^{\mathbf{T}'}(u), \ell^{\mathbf{T}'}(u \cdot 2)) \in K$, we have $\ell^{\mathbf{T}'}(u \cdot 2) = (c_1 \dots c_m, 1)$ for some $c_1 \dots c_m \in F$. Since for all $i \leq n-1$ we must have $(\ell^{\mathbf{T}'}(u \cdot 2 \cdot 1^{i-1}), \ell^{\mathbf{T}'}(u \cdot 2 \cdot 1^i)) \in J$, we get $\ell^{\mathbf{T}'}(u \cdot 2 \cdot 1^{i-1}) = (c_1 \dots c_m, i) \in \Sigma'$ for $i = 1, \dots, n$. This implies $n \leq m$. But $\ell^{\mathbf{T}'}(u \cdot 2 \cdot 1^{n-1}) = (c_1 \dots c_m, n)$ must be in Y , so $m = n$. Since $\pi((c_1 \dots c_n, i)) = c_i = \ell^{\mathbf{T}}(u \cdot 2 \cdot 1^{i-1})$, it follows that (2) holds with \mathbf{T}' in place of $\hat{\mathbf{T}}$. \square

We assume the standard definition of the *yield* function $\mathbf{y}: \mathbb{T}_\Sigma \rightarrow \Sigma^+$. Using the term notation for unranked trees, we can define it as follows:

$$\begin{aligned} \mathbf{y}(c) &= c, \\ \mathbf{y}(c(t_1 \dots t_n)) &= \mathbf{y}(t_1) \dots \mathbf{y}(t_n). \end{aligned}$$

As is well known, a set of non-empty strings is a context-free language if and only if it is the yield image of a local set of trees.

Let us call a tree $\mathbf{T} = (T, \ell) \in \mathbb{T}_\Sigma$ *disjointly labeled with* Σ_0, Σ_1 if (i) Σ_0 and Σ_1 are disjoint subsets of Σ , (ii) $u \in \text{dom}(\prec_2^T)$ implies $\ell(u) \in \Sigma_1$, and (iii) $u \in T - \text{dom}(\prec_2^T)$ implies $\ell(u) \in \Sigma_0$. Let Σ_0, Σ_1 be disjoint sets and let

$$\mathbb{T}_{\Sigma_0}^{\Sigma_1} = \{ \mathbf{T} \in \mathbb{T}_{\Sigma_0 \cup \Sigma_1} \mid \mathbf{T} \text{ is disjointly labeled with } \Sigma_0, \Sigma_1 \}.$$

On $\mathbb{T}_{\Sigma_0}^{\Sigma_1}$, the yield function $\mathbf{y}: \mathbb{T}_{\Sigma_0}^{\Sigma_1} \rightarrow \Sigma_0^+$ can be expressed as the composition¹¹

$$\mathbf{y} = \mathbf{h}_{\Sigma_0, \Sigma_1} \circ \mathbf{enc}$$

of the string encoding function \mathbf{enc} and an alphabetic homomorphism $\mathbf{h}_{\Sigma_0, \Sigma_1}: (\Gamma_{\Sigma_0 \cup \Sigma_1})^* \rightarrow \Sigma_0^*$ defined by

$$\mathbf{h}_{\Sigma_0, \Sigma_1}(\mathbb{1}_c) = \begin{cases} c & \text{if } c \in \Sigma_0, \\ \varepsilon & \text{if } c \in \Sigma_1, \end{cases}$$

$$\mathbf{h}_{\Sigma_0, \Sigma_1}(\mathbb{1}_c) = \varepsilon.$$

Lemma 4. *Let $L \subseteq \Sigma^+$ be a context-free language. There exist a set Υ disjoint from Σ and a local set $K \subseteq \mathbb{T}_\Sigma^\Upsilon$ such that $L = \mathbf{y}(K) = \mathbf{h}_{\Sigma, \Upsilon}(\mathbf{enc}(K))$.*

Proof. Let $G = (N, \Sigma, P, S)$ be an ε -free context-free grammar for L . Clearly, the parse trees of G form a local subset K of \mathbb{T}_Σ^N , and $L = \mathbf{y}(K) = \mathbf{h}_{\Sigma, N}(\mathbf{enc}(K))$.¹² \square

Conversely, $h(\mathbf{enc}(K))$ is always a context-free language whenever K is a local set of trees and h is a homomorphism.

A projection $\pi: \Sigma' \rightarrow \Sigma$ induces a projection $\hat{\pi}: \Gamma_{\Sigma'} \rightarrow \Gamma_\Sigma$ in an obvious way:

$$\hat{\pi}(\mathbb{1}_c) = \mathbb{1}_{\pi(c)}, \quad \hat{\pi}(\mathbb{1}_c) = \mathbb{1}_{\pi(c)}.$$

Lemma 5. *Let $\pi: \Sigma' \rightarrow \Sigma$ be a projection and $L \subseteq \mathbb{T}_{\Sigma'}$. Then $\mathbf{enc}(\pi(L)) = \hat{\pi}(\mathbf{enc}(L))$.*

We can use Lemmas 2 through 5 to show that every context-free language L can be represented as $L = h(R \cap D'_n)$ with some alphabetic homomorphism h and local set R . The Chomsky-Schützenberger theorem, however, is stated in terms of the Dyck language D_n rather than the set D'_n of Dyck primes. The following easy lemma bridges the representation in terms of D'_n and that in terms of D_n .

¹¹ If $f: A \rightarrow B$ and $g: B \rightarrow C$ are functions, I write $g \circ f$ for the composition of f and g defined by $(g \circ f)(x) = g(f(x))$. Note that the order of the two functions is reversed compared to the case of compositions of binary relations.

¹² This also follows from the fact that a local set of trees is always obtained from a local set of disjointly labeled trees by a projection that does not change the labels of leaves.

Lemma 6. *Let $L \subseteq \Gamma_{\Sigma}^+$ be a local set of strings. Then there exist a finite alphabet Σ' , a projection $\pi: \Sigma' \rightarrow \Sigma$, and a local set $L' \subseteq \Gamma_{\Sigma'}^+$ such that $L \cap D'_{\Sigma} = \widehat{\pi}(L' \cap D_{\Sigma'})$. Moreover, $\widehat{\pi}$ maps $L' \cap D_{\Sigma'}$ bijectively to $L \cap D'_{\Sigma}$.*

Proof. Let $A, Z \subseteq \Gamma_{\Sigma}, I \subseteq \Gamma_{\Sigma}^2$ be finite sets such that $L = A\Gamma_{\Sigma}^* \cap \Gamma_{\Sigma}^*Z - (\Gamma_{\Sigma}^*(\Gamma_{\Sigma}^2 - I)\Gamma_{\Sigma}^*)$. We may assume without loss of generality that Σ is finite. Let $\Sigma' = \Sigma \cup \{\bar{c} \mid c \in \Sigma\}$. Let

$$\begin{aligned} A' &= \{ \llbracket \bar{c} \mid \llbracket c \in A \}, \\ Z' &= \{ \rrbracket \bar{c} \mid \rrbracket c \in Z \}, \\ I' &= I \cup \{ \llbracket \bar{c} d \mid \llbracket c d \in I \} \cup \{ d \rrbracket \bar{c} \mid d \rrbracket c \in I \} \cup \{ \llbracket \bar{c} \rrbracket \bar{c} \mid \llbracket c \rrbracket c \in I \}, \end{aligned}$$

and put

$$L' = A'\Gamma_{\Sigma'}^* \cap \Gamma_{\Sigma'}^*Z' - (\Gamma_{\Sigma'}^*(\Gamma_{\Sigma'}^2 - I')\Gamma_{\Sigma'}^*).$$

Define $\pi: \Sigma' \rightarrow \Sigma$ by

$$\pi(c) = c, \quad \pi(\bar{c}) = c$$

for each $c \in \Sigma$. It is easy to check that L' and π satisfy the desired properties. \square

Lemma 7. *If $L \subseteq \mathbb{T}_{\Sigma}$ is a local set, then there exist a finite alphabet Σ' , a projection $\pi: \Sigma' \rightarrow \Sigma$, and a local set $L' \subseteq \Gamma_{\Sigma'}^+$ such that $\mathbf{enc}(L) = \widehat{\pi}(L' \cap D_{\Sigma'})$. Moreover, $\mathbf{enc}^{-1} \circ \widehat{\pi}$ maps $L' \cap D_{\Sigma'}$ bijectively to L .*

Proof. By Lemma 3, there exist a projection $\pi_1: \Sigma_1 \rightarrow \Sigma$ and a super-local set $L_1 \subseteq \mathbb{T}_{\Sigma_1}$ such that $L = \pi_1(L_1)$ and π_1 is a bijection from L_1 to L . By Lemma 2, there exists a local set $L_2 \subseteq \Gamma_{\Sigma_1}^*$ such that $\mathbf{enc}(L_1) = L_2 \cap D'_{\Sigma_1}$. By Lemma 6, there exist a projection $\pi_3: \Sigma' \rightarrow \Sigma_1$ and a local set $L' \subseteq \Gamma_{\Sigma'}^*$ such that $L_2 \cap D'_{\Sigma_1} = \widehat{\pi_3}(L' \cap D_{\Sigma'})$ and $\widehat{\pi_3}$ is a bijection from $L' \cap D_{\Sigma'}$ to $L_2 \cap D'_{\Sigma_1}$. By Lemma 5, $\mathbf{enc}(L) = \mathbf{enc}(\pi_1(L_1)) = \widehat{\pi_1}(\mathbf{enc}(L_1))$. Since \mathbf{enc} is injective, $\widehat{\pi_1}$ is a bijection from $\mathbf{enc}(L_1)$ to $\mathbf{enc}(L)$. Taking these all together, we get

$$\begin{aligned} \mathbf{enc}(L) &= \mathbf{enc}(\pi_1(L_1)) \\ &= \widehat{\pi_1}(\mathbf{enc}(L_1)) \\ &= \widehat{\pi_1}(L_2 \cap D'_{\Sigma_1}) \\ &= \widehat{\pi_1}(\widehat{\pi_3}(L' \cap D_{\Sigma'})) \\ &= (\widehat{\pi_1} \circ \widehat{\pi_3})(L' \cap D_{\Sigma'}) \\ &= \widehat{\pi}(L' \cap D_{\Sigma'}), \end{aligned}$$

where $\pi = \pi_1 \circ \pi_3$. Since $\widehat{\pi_3}$ is a bijection from $L' \cap D_{\Sigma'}$ to $L_2 \cap D'_{\Sigma_1} = \mathbf{enc}(L_1)$ and $\widehat{\pi_1}$ is a bijection from $\mathbf{enc}(L_1)$ to $\mathbf{enc}(L)$, $\widehat{\pi}$ is a bijection from $L' \cap D_{\Sigma'}$ to $\mathbf{enc}(L)$, and the second statement of the lemma follows. \square

Theorem 8 (Chomsky and Schützenberger). *A language $L \subseteq \Sigma^+$ is context-free if and only if there exist a positive integer n , a local set $R \subseteq \Gamma_n^+$, and an alphabetic homomorphism $h: \Gamma_n^* \rightarrow \Sigma^*$ such that $L = h(R \cap D_n)$.*

Proof. The “if” direction is by standard closure properties of the context-free languages. For the “only if” direction, let $L \subseteq \Sigma^+$ be a context-free language. Then Lemma 4 gives an alphabet \mathcal{T} disjoint from Σ and a local set $K \subseteq \mathbb{T}_{\Sigma}^{\mathcal{T}}$ such that $L = \mathbf{h}_{\Sigma, \mathcal{T}}(\mathbf{enc}(K))$. By Lemma 7, there are a projection $\pi: \mathcal{Y}' \rightarrow \Sigma \cup \mathcal{Y}$ and a local set $R \subseteq \Gamma_{\mathcal{Y}'}^+$ such that $\mathbf{enc}(K) = \widehat{\pi}(R \cap D_{\mathcal{Y}'})$. We have

$$\begin{aligned} L &= \mathbf{h}_{\Sigma, \mathcal{T}}(\mathbf{enc}(K)) \\ &= \mathbf{h}_{\Sigma, \mathcal{T}}(\widehat{\pi}(R \cap D_{\mathcal{Y}'})), \end{aligned}$$

so the required condition holds with $n = |\mathcal{Y}'|$ and $h = \mathbf{h}_{\Sigma, \mathcal{T}} \circ \widehat{\pi}$.¹³ \square

In the proof of Theorem 8, $\mathbf{enc}^{-1} \circ \widehat{\pi}$ is a bijection from $R \cap D_{\mathcal{Y}'}$ to K . (See the second statement in Lemma 7.) Given how $\widehat{\pi}$ is defined in this proof, it is clear that the elements of $R \cap D_{\mathcal{Y}'}$ are simply the elements of K expressed in alternative notation. If K is the set of derivation trees of a context-free grammar for L , then an element s of $R \cap D_{\mathcal{Y}'}$ represents both the element $t = \mathbf{enc}^{-1}(\widehat{\pi}(s))$ of K and the element $h_{\Sigma, \mathcal{T}}(\widehat{\pi}(s)) = h_{\Sigma, \mathcal{T}}(\mathbf{enc}(t)) = \mathbf{y}(t)$ of L . Moreover, every pair $(t, \mathbf{y}(t))$ with $t \in K$ is so represented. This is an important consequence of the theorem explicitly noted by Chomsky [5, page 377], though rarely emphasized since.¹⁴

We took a rather long route to the Chomsky-Schützenberger theorem. Our generalization of the theorem to multi-dimensional tree languages follows a similar path, except that an analogue of Lemma 4 is not needed, since the multi-dimensional counterpart of the function \mathbf{enc} is not exactly a generalization of the usual notion.

4 Multi-dimensional Trees

Multi-dimensional trees were introduced by Baldwin and Strawn [2] and further investigated by Rogers [33,32] in connection with tree-adjointing grammars. In an ordinary (labeled, ordered unranked) tree, the set of children of a node forms a linearly ordered sequence of labeled nodes, i.e., a string. In an m -dimensional tree ($m \geq 1$), the set of children of a node (if non-empty) forms an $(m - 1)$ -dimensional tree. A 0-dimensional tree just consists of a single labeled node.

Unlike Rogers [33,32], who introduces the higher-dimensional tree as a new kind of object, I prefer to define an m -dimensional tree as a special kind of ordinary m -ary tree.¹⁵ This corresponds to one of the encodings of m -dimensional

¹³ Here, $|\mathcal{Y}'|$ denotes the cardinality of the set \mathcal{Y}' . See footnote 10.

¹⁴ Instead of a super-local set of trees, Chomsky [5] used the notion of a *modified normal grammar*, a restricted kind of grammar in Chomsky normal form.

¹⁵ To be precise, our m -dimensional trees form a special class of m -ary *cardinal trees* in the sense of Benoit et al. [3]. In m -ary cardinal trees, each node has m slots for children, each of which may or may not be occupied, independently of the other slots. Cardinal trees are also known as *tries*.

trees considered by Baldwin and Strawn. The first-child-next-sibling encoding of unranked trees will be a special case of this definition for $m = 2$.¹⁶

We use finite strings of elements of $[1, m] = \{1, \dots, m\}$ to represent nodes of m -ary trees. We write $u \cdot v$ for the concatenation of finite strings u, v over $[1, m]$, and write ε for the empty string.

An m -ary tree domain is any non-empty, finite, prefix-closed subset of $[1, m]^*$. (Since $\emptyset^* = \{\varepsilon\}$, the only 0-ary tree domain is $\{\varepsilon\}$.) If T is an m -ary tree domain, we write $u \prec_i^T v$ to mean $u, v \in T$ and $u \cdot i = v$. If Σ is a (possibly infinite) set of symbols, an m -ary tree over Σ is a pair (T, ℓ) , where T is an m -ary tree domain and ℓ is a function from T to Σ .

If $\mathbf{T} = (T, \ell^{\mathbf{T}})$ is an m -ary tree and $U \subseteq T$ is an m -ary tree domain, then the *restriction of \mathbf{T} to U* is the m -ary tree

$$\mathbf{T} \upharpoonright U = (U, \ell^{\mathbf{T}} \upharpoonright U).$$

When $u \in T$, let

$$T/u = \{v \mid uv \in T\}.$$

Then T/u is an m -ary tree domain and the *subtree of \mathbf{T} rooted at u* is defined by

$$\mathbf{T}/u = (T/u, \ell),$$

where $\ell(v) = \ell^{\mathbf{T}}(uv)$.

Recall that a first-child-next-sibling encoding of an unranked tree is a binary tree (T, ℓ) such that $1 \notin T$. Analogously, an m -dimensional tree is an m -ary tree (T, ℓ) such that T is an m -ary tree domain included in a certain special subset of $[1, m]^*$. For each natural number m , define a subset \mathbb{P}_m of $[1, m]^*$ by induction, as follows:

$$\begin{aligned} \mathbb{P}_0 &= \{\varepsilon\}, \\ \mathbb{P}_m &= (m \cdot \mathbb{P}_{m-1})^* \quad \text{for } m \geq 1. \end{aligned}$$

It is easy to see that $w \in \mathbb{P}_m$ if and only if $w = i \cdot v$ implies $i = m$ and $w = u \cdot i \cdot j \cdot v$ implies $j \geq i - 1$. For $m \geq 0$, an m -dimensional tree (over Σ) is an m -ary tree (over Σ) $\mathbf{T} = (T, \ell^{\mathbf{T}})$ such that $T \subseteq \mathbb{P}_m$. For $m \geq 1$, an m -dimensional hedge (over Σ) is an m -ary tree (over Σ) $\mathbf{T} = (T, \ell^{\mathbf{T}})$ such that $T \subseteq \mathbb{P}_{m-1} \cdot \mathbb{P}_m$.¹⁷ We write \mathbb{T}_Σ^m and \mathbb{H}_Σ^m to denote the set of all m -dimensional trees over Σ and the set of all m -dimensional hedges over Σ , respectively. For $m \geq 1$, all m -dimensional trees are m -dimensional hedges. Note that a 1-dimensional hedge is just a 1-dimensional tree.

Note that a 0-dimensional tree is a structure $\mathbf{T}_c = (\{\varepsilon\}, \{(\varepsilon, c)\})$ consisting of a single node labeled by some $c \in \Sigma$. We may identify \mathbf{T}_c with c ; under this convention, $\mathbb{T}_\Sigma^0 = \Sigma$. Note that if $m \neq n$, then $\mathbb{T}_\Sigma^m \cap \mathbb{T}_\Sigma^n = \mathbb{T}_\Sigma^0$.

A 1-dimensional tree is a non-empty, linearly ordered sequence of labeled nodes. We may use a string $c_1 \dots c_n \in \Sigma^+$ to denote the 1-dimensional tree

¹⁶ This point was already noted by Kasprzik [20].

¹⁷ Our m -dimensional hedges correspond to Baldwin and Strawn's [2] multidimensional forests of dimension m and degree $m - 1$.

$\mathbf{T}_{c_1 \dots c_n} = (\{\varepsilon, 1, \dots, 1^{n-1}\}, \ell)$, where $\ell(1^{i-1}) = c_i$ for $i = 1, \dots, n$. Under this convention, $\mathbb{T}_\Sigma^1 = \Sigma^+$.

The first-child-next-sibling encodings of unranked trees and unranked hedges coincide with the 2-dimensional trees and the 2-dimensional hedges, respectively; we have $\mathbb{T}_\Sigma^2 = \mathbb{T}_\Sigma$.

Henceforth, we use $\mathbf{T}, \mathbf{T}', \mathbf{U}$, etc., as variables ranging over m -dimensional trees and m -dimensional hedges. Unless we indicate otherwise, we assume $\mathbf{T} = (T, \ell^{\mathbf{T}})$, $\mathbf{T}' = (T', \ell^{\mathbf{T}'})$, $\mathbf{U} = (U, \ell^{\mathbf{U}})$, etc.

Let \mathbf{T} be an m -dimensional hedge. We can see that if $u \in T$, then

$$ST_i(\mathbf{T}, u) = (\mathbf{T}/u) \upharpoonright \{v \in \mathbb{P}_i \mid uv \in T\}$$

is always an i -dimensional tree, and

$$SH_i(\mathbf{T}, u) = (\mathbf{T}/u) \upharpoonright \{v \in \mathbb{P}_{i-1} \cdot \mathbb{P}_i \mid uv \in T\}$$

is always an i -dimensional hedge. In particular, when $u \cdot m \in T$, the subtree $\mathbf{T}/(u \cdot m) = SH_m(\mathbf{T}, u \cdot m)$ is always an m -dimensional hedge.

For $i \geq 1$, we write $u \triangleleft_i^T v$ to mean

$$v \in T \cap u \cdot i \cdot \mathbb{P}_{i-1}.$$

When $u \triangleleft_i^T v$, we say that v is a *child of u in the i -th dimension* (in \mathbf{T}). If $u \prec_i^T v$, then v is the *first child of u in the i -th dimension*. Define

$$C_i^T(u) = \{v \in \mathbb{P}_{i-1} \mid u \cdot i \cdot v \in T\} = \{v \mid u \triangleleft_i^T u \cdot i \cdot v\}.$$

If $u \cdot i \notin T$, that is, if $u \notin \text{dom}(\prec_i^T)$, then $C_i^T(u) = \emptyset$. If $u \cdot i \in T$, define

$$\mathbf{C}_i^T(u) = ST_{i-1}(\mathbf{T}, u \cdot i) = \mathbf{T}/(u \cdot i) \upharpoonright C_i^T(u).$$

Then $\mathbf{C}_i^T(u)$ is always an $(i-1)$ -dimensional tree.

We assume that elements of $[1, m]^*$ are alphabetically ordered, with $k+1$ alphabetically *preceding* k . We write $u \triangleleft_{i,j}^T v$ to mean v is the j -th node, under this ordering, among the children of u in the i -th dimension. The *degree* of a node $v \in T$ is the number of children of v in the m -th (i.e., highest) dimension.

A subset of \mathbb{T}_Σ^m is an *m -dimensional tree language*. We allow Σ to be an infinite set, but are usually interested in m -dimensional tree languages over some finite subset of Σ .

We call a set $L \subseteq \mathbb{T}_\Sigma^m$ *degree-bounded* if there exists a k such that for all $\mathbf{T} \in L$ and for all $v \in T$, the degree of v does not exceed k .

It is sometimes helpful to use term-like notations for m -dimensional hedges and trees. Let P be an $(m-1)$ -ary tree domain included in \mathbb{P}_{m-1} (i.e., a finite, non-empty, prefix-closed subset of \mathbb{P}_{m-1}), and let u_1, \dots, u_k be the elements of P , in alphabetical order (which implies $u_1 = \varepsilon$). If $\mathbf{T}_1, \dots, \mathbf{T}_k \in \mathbb{T}_\Sigma^m$, then we write

$$P(\mathbf{T}_1, \dots, \mathbf{T}_k)$$

to denote the m -dimensional hedge $\mathbf{U} = (U, \ell^{\mathbf{U}}) \in \mathbb{H}_{\Sigma}^m$ such that

$$U = \bigcup_{i=1}^k u_i \cdot T_i, \quad ST_m(\mathbf{U}, u_i) = T_i.$$

(As a degenerate case, we have $\{\varepsilon\}(\mathbf{T}) = \mathbf{T}$ for any $\mathbf{T} \in \mathbb{T}_{\Sigma}^m$.) If $\mathbf{T} \in \mathbb{H}_{\Sigma}^m$ and $c \in \Sigma$, then we write

$$c -_m \mathbf{T}$$

to denote the m -dimensional tree $\mathbf{V} = (V, \ell^{\mathbf{V}}) \in \mathbb{T}_{\Sigma}^m$ such that

$$V = \{\varepsilon\} \cup m \cdot T, \quad \ell^{\mathbf{V}}(\varepsilon) = c, \quad SH_m(\mathbf{V}, m) = \mathbf{T}.$$

Combining the two notations,

$$c -_m P(\mathbf{T}_1, \dots, \mathbf{T}_k)$$

denotes the m -dimensional tree $\mathbf{T} = (T, \ell^{\mathbf{T}})$ such that $\ell^{\mathbf{T}}(\varepsilon) = c$, $C_m^T(\varepsilon) = P$, and $ST_m(\mathbf{T}, u_i) = T_i$, where u_1, \dots, u_k is the alphabetical listing of the elements of P .¹⁸

Example 9. Derivation trees of a simple context-free tree grammar $G = (N, \Sigma, P, S)$ can be represented as 3-dimensional trees over the alphabet $N \cup \Sigma$. In these 3-dimensional trees, a node has children in the third dimension if and only if it is labeled by a nonterminal. For instance, the derivation tree $\pi_1(\pi_2(\pi_3))$ of the grammar from Example 1 may be represented by the 3-dimensional tree \mathbf{T} in Fig. 2. In this tree, the node labeled by S is the root; the edges in the third dimension are drawn as dotted lines, those in the second dimension solid, and those in the first dimension dashed. For instance, the node 3 (i.e., the first child of the root in the third dimension) is labeled by the nonterminal B , and its children in the third dimension form a 2-dimensional tree corresponding to the right-hand side of the rule $\pi_2 = B(x_1x_2) \rightarrow h(a_1B(h(a_2x_1a_3)h(a_4x_2a_5))a_6)$. The numbering of variables in the rules are eschewed in favor of a single variable \mathbf{x} ; the alphabetic ordering of the nodes labeled by \mathbf{x} among the children in the third dimension of a nonterminal-labeled node is assumed to correspond to the numbering.¹⁹ This tree can be represented in the term notation as follows (omitting the dots in the strings over $\{1, 2\}$ representing nodes):

$$\begin{aligned} & S -_3 \{\varepsilon, 2, 21\}(\\ & \quad B -_3 \{\varepsilon, 2, 21, 212, 2122, 21221, 212211, 2121, 21212, 212121, 2121211, 211\}(\\ & \quad \quad h, a_1, B -_3 \{\varepsilon, 2, 21\}(g, \mathbf{x}, \mathbf{x}), h, a_2, \mathbf{x}, a_3, h, a_4, \mathbf{x}, a_5, a_6 \\ & \quad \quad \quad), \\ & \quad \quad \quad e, \\ & \quad \quad \quad e \\ & \quad \quad \quad). \end{aligned}$$

¹⁸ An equivalent notation for m -dimensional trees has been used by Kasprzik [21].

¹⁹ It is known that simple context-free tree grammars satisfying this condition constitute a normal form. This use of a single variable instead of numbered variables is not crucial for our purposes.

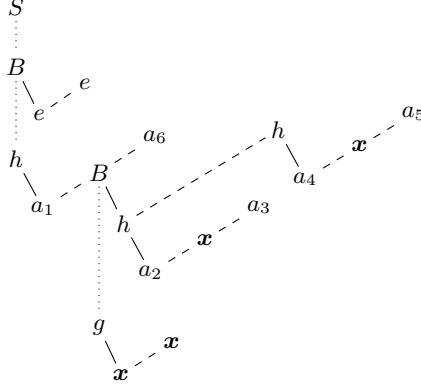


Fig. 2. A derivation tree of a simple context-free tree grammar represented as a 3-dimensional tree.

We have

$$\begin{aligned} \mathbf{C}_3^T(3 \cdot 3 \cdot 2 \cdot 1) &= g_{-2} \{\varepsilon, 1\}(\mathbf{x}, \mathbf{x}) \\ &= g(\mathbf{x}\mathbf{x}), \\ \mathbf{C}_2^T(3 \cdot 3 \cdot 2 \cdot 1) &= h_{-1} h \\ &= hh, \end{aligned}$$

using both the notation introduced just above and the standard term and string representations for 2-dimensional and 1-dimensional trees.

5 Local and Super-local Sets of Multi-dimensional Trees

If $A, Z \subseteq \Sigma$ and $I \subseteq \Sigma \times \mathbb{T}_\Sigma^{m-1}$ are finite sets, we let $\text{Loc}^m(A, Z, I)$ denote the set of all m -dimensional trees $\mathbf{T} = (T, \ell^{\mathbf{T}})$ in \mathbb{T}_Σ^m that satisfy the following conditions:

- L1. $\ell^{\mathbf{T}}(\varepsilon) \in A$,
- L2. $v \in T - \text{dom}(\prec_m^{\mathbf{T}})$ implies $\ell^{\mathbf{T}}(v) \in Z$, and
- L3. $v \in \text{dom}(\prec_m^{\mathbf{T}})$ implies $(\ell^{\mathbf{T}}(v), \mathbf{C}_m^{\mathbf{T}}(v)) \in I$.

A set $L \subseteq \mathbb{T}_\Sigma^m$ is *local* [33,32] if there exist finite sets $A, Z \subseteq \Sigma$ and $I \subseteq \Sigma \times \mathbb{T}_\Sigma^{m-1}$ such that $L = \text{Loc}^m(A, Z, I)$. Note that if $L \subseteq \mathbb{T}_\Sigma^m$ is local, then L must be degree-bounded; for, if $L = \text{Loc}^m(A, Z, I)$, the degree of any node v of $\mathbf{T} \in L$ is bounded by the maximal size of \mathbf{U} such that $(c, \mathbf{U}) \in I$ for some c . Clearly, the notion of locality coincides with the usual notion [31,42,40] when $m \in \{1, 2\}$.

Let $m \geq 2$. Write \mathbb{N}_+ for $\mathbb{N} - \{0\}$ (the set of positive integers). If $A, Z, Y \subseteq \Sigma$, $K \subseteq \Sigma \times \Sigma$, and $J \subseteq \Sigma \times \{P \subseteq \mathbb{P}_{m-2} \mid P \text{ is an } (m-2)\text{-ary tree domain}\} \times \mathbb{N}_+ \times \Sigma$ are finite sets, then we let $\text{SLoc}^m(A, Z, K, Y, J)$ denote the set of all trees $\mathbf{T} = (T, \ell^{\mathbf{T}})$ in \mathbb{T}_Σ^m that satisfy the following conditions:

- SL1. $\ell^{\mathbf{T}}(\varepsilon) \in A$,
 SL2. $v \in T - \text{dom}(\prec_m^T)$ implies $\ell^{\mathbf{T}}(v) \in Z$,
 SL3. $u \prec_m^T v$ implies $(\ell^{\mathbf{T}}(u), \ell^{\mathbf{T}}(v)) \in K$,
 SL4. $v \neq \varepsilon$ and $v \in T - \text{dom}(\prec_{m-1}^T)$ imply $\ell^{\mathbf{T}}(v) \in Y$, and
 SL5. $u \in \text{dom}(\prec_{m-1}^T)$ and $u \prec_{m-1,i}^T v$ imply $(\ell^{\mathbf{T}}(u), C_{m-1}^T(u), i, \ell^{\mathbf{T}}(v)) \in J$.

We call a set $L \subseteq \mathbb{T}_{\Sigma}^m$ *super-local* if there exist finite sets $A, Z, Y \subseteq \Sigma$, $K \subseteq \Sigma \times \Sigma$, and $J \subseteq \Sigma \times \{P \subseteq \mathbb{P}_{m-2} \mid P \text{ is an } (m-2)\text{-ary tree domain}\} \times \mathbb{N}_+ \times \Sigma$ such that $L = \text{SLoc}^m(A, Z, K, Y, J)$. For $m = 2$, $\mathbb{P}_{m-2} = \mathbb{P}_0 = \{\varepsilon\}$, and $u \prec_{1,i}^T v$ only if $i = 1$, so this definition generalizes our earlier definition of super-locality for subsets of $\mathbb{T}_{\Sigma} = \mathbb{T}_{\Sigma}^2$. It is easy to see that a degree-bounded super-local language must be local. Although we allow Σ to be infinite, any local or super-local set $L \subseteq \mathbb{T}_{\Sigma}^m$ must be an m -dimensional tree language over some finite subset of Σ .

Projections from Σ to Σ' are naturally extended to m -dimensional trees and hedges over Σ and to m -dimensional tree languages over Σ . The next lemma generalizes Lemma 3 to the higher-dimensional case. The proofs of two lemmas to follow (Lemmas 11 and 32) will be adaptations of the proof of this lemma.

Lemma 10. *Let $m \geq 2$. For any local m -dimensional tree language $L \subseteq \mathbb{T}_{\Sigma}^m$, there exist a finite alphabet Σ' , a degree-bounded, super-local m -dimensional tree language $L' \subseteq \mathbb{T}_{\Sigma'}^m$, and a projection $\pi: \Sigma' \rightarrow \Sigma$ such that $L = \pi(L')$. Moreover, π maps L' bijectively to L .*

Proof. The proof parallels that of Lemma 3. The idea is to change the label of each non-root node v of $\mathbf{T} \in L$ to

$$(C_m^{\mathbf{T}}(u), v'),$$

where $u \cdot m \cdot v' = v$ and $v' \in \mathbb{P}_{m-1}$. For uniformity, we change the label of the root from $c \in \Sigma$ to $(\mathbf{T}_c, \varepsilon)$, where $\mathbf{T}_c = (\{\varepsilon\}, \{(\varepsilon, c)\})$ is the single-node tree that we identified with c . The relabeled m -dimensional trees obtained this way form a super-local set, and we can get back the original m -dimensional trees by a projection.

Let

$$\Sigma'' = \{(\mathbf{T}, v) \mid \mathbf{T} = (T, \ell^{\mathbf{T}}) \in \mathbb{T}_{\Sigma}^{m-1}, v \in T\},$$

and define a projection $\pi: \Sigma'' \rightarrow \Sigma$ by

$$\pi((\mathbf{T}, v)) = \ell^{\mathbf{T}}(v).$$

Suppose that $A, Z \subseteq \Sigma$ and $I \subseteq \Sigma \times \mathbb{T}_{\Sigma}^{m-1}$ are finite sets such that $L = \text{Loc}^m(A, Z, I)$. Let

$$F = \{\mathbf{T}_c \mid c \in A\} \cup \{\mathbf{T} \mid (c, \mathbf{T}) \in I\},$$

$$\Sigma' = \{(\mathbf{T}, v) \in \Sigma'' \mid \mathbf{T} \in F\}.$$

Note that Σ' is a finite subset of Σ'' . Now define

$$A' = \{(\mathbf{T}_c, \varepsilon) \mid c \in A\},$$

$$\begin{aligned}
Z' &= \{ (\mathbf{T}, v) \in \Sigma' \mid \ell^{\mathbf{T}}(v) \in Z \}, \\
K &= \{ ((\mathbf{T}, v), (\mathbf{T}', \varepsilon)) \mid (\mathbf{T}, v) \in \Sigma', (\ell^{\mathbf{T}}(v), \mathbf{T}') \in I \}, \\
Y &= \{ (\mathbf{T}, v) \in \Sigma' \mid v \notin \text{dom}(\prec_{m-1}^{\mathbf{T}}) \}, \\
J &= \{ ((\mathbf{T}, u), C_{m-1}^{\mathbf{T}}(u), i, (\mathbf{T}, v)) \mid (\mathbf{T}, u) \in \Sigma', u \prec_{m-1, i}^{\mathbf{T}} v \}.
\end{aligned}$$

These are all finite sets. Let $L' \subseteq \mathbb{T}_{\Sigma'}^m$, be the super-local set defined by $L' = \text{SLoc}^m(A', Z', K, Y, J)$. It is quite clear that $L' \subseteq \mathbb{T}_{\Sigma'}^m$. We show that L' and π (restricted to Σ') satisfy the required properties.

For each $\mathbf{T} \in \mathbb{T}_{\Sigma}^m$, define an m -dimensional tree $\hat{\mathbf{T}} = (T, \ell^{\hat{\mathbf{T}}}) \in \mathbb{T}_{\Sigma'}^m$ by

$$\ell^{\hat{\mathbf{T}}}(\varepsilon) = (\mathbf{T}_{\ell^{\mathbf{T}}(\varepsilon)}, \varepsilon), \quad (4)$$

$$\ell^{\hat{\mathbf{T}}}(u \cdot m \cdot v) = (C_m^{\mathbf{T}}(u), v), \quad \text{if } u \in \text{dom}(\prec_m^{\mathbf{T}}) \text{ and } v \in C_{m-1}^{\mathbf{T}}(u). \quad (5)$$

It is clear that $\pi(\hat{\mathbf{T}}) = \mathbf{T}$ for all $\mathbf{T} \in \mathbb{T}_{\Sigma}^m$. Our goal is to show

$$L' = \{ \hat{\mathbf{T}} \mid \mathbf{T} \in L \}.$$

This clearly implies that π is a bijection from L' to L .

We show that for all $\mathbf{T} \in \mathbb{T}_{\Sigma}^m$,

$$\mathbf{T} \in L \text{ if and only if } \hat{\mathbf{T}} \in L'. \quad (6)$$

This follows from five observations. Firstly, note the following:

- Suppose $u \in T - \text{dom}(\prec_{m-1}^{\hat{\mathbf{T}}})$. If $\ell^{\hat{\mathbf{T}}}(u) = (\mathbf{U}, v)$, then $v \notin \text{dom}(\prec_{m-1}^{\mathbf{U}})$. This means that $\ell^{\hat{\mathbf{T}}}(u) \in Y$ if $\ell^{\hat{\mathbf{T}}}(u) \in \Sigma'$.
- Suppose $s \prec_m^{\hat{\mathbf{T}}} u = s \cdot m \cdot u' \prec_{m-1, i}^{\hat{\mathbf{T}}} v = u \cdot (m-1) \cdot v'$. Then we have $u' \prec_{m-1, i}^{C_m^{\mathbf{T}}(s)} u' \cdot (m-1) \cdot v'$ and

$$\begin{aligned}
\ell^{\hat{\mathbf{T}}}(u) &= (C_m^{\mathbf{T}}(s), u'), \\
\ell^{\hat{\mathbf{T}}}(v) &= (C_m^{\mathbf{T}}(s), u' \cdot (m-1) \cdot v'), \\
C_{m-1}^{\hat{\mathbf{T}}}(u) &= C_{m-1}^{C_m^{\mathbf{T}}(s)}(u').
\end{aligned}$$

This means that $(\ell^{\hat{\mathbf{T}}}(u), C_{m-1}^{\hat{\mathbf{T}}}(u), i, \ell^{\hat{\mathbf{T}}}(v)) \in J$ if $\ell^{\hat{\mathbf{T}}}(u) \in \Sigma'$.

Thus, $\hat{\mathbf{T}}$ satisfies the last two conditions SL4 and SL5 for membership in $\text{SLoc}^m(A', Z', K, Y, J)$ whenever $\hat{\mathbf{T}} \in \mathbb{T}_{\Sigma'}^m$. Secondly, the following biconditional always holds:

- $\ell^{\hat{\mathbf{T}}}(\varepsilon) \in A$ if and only if $\ell^{\hat{\mathbf{T}}}(\varepsilon) \in A'$.

Thirdly, the following biconditional holds whenever $\ell^{\hat{\mathbf{T}}}(v) \in \Sigma'$:

- $\ell^{\hat{\mathbf{T}}}(v) \in Z$ if and only if $\ell^{\hat{\mathbf{T}}}(v) \in Z'$.

Fourthly, if $u \prec_m^T v$ and $\ell^{\hat{\mathbf{T}}}(u) \in \Sigma'$, then the following biconditional holds:

$$- (\ell^{\mathbf{T}}(u), \mathbf{C}_m^{\mathbf{T}}(u)) \in I \text{ if and only if } (\ell^{\hat{\mathbf{T}}}(u), \ell^{\hat{\mathbf{T}}}(v)) \in K.$$

Lastly, it is easy to see that $\mathbf{T} \in L$ implies $\hat{\mathbf{T}} \in \mathbb{T}_{\Sigma'}^m$. Combining these five observations, we get (6).

It follows from the ‘‘only if’’ direction of (6) that $\{\hat{\mathbf{T}} \mid \mathbf{T} \in L\} \subseteq L'$. To establish the converse inclusion, we show that

$$\text{if } \mathbf{T}' \in L' \text{ and } \mathbf{T} = \pi(\mathbf{T}'), \text{ then } \mathbf{T}' = \hat{\mathbf{T}}.$$

This together with the ‘‘if’’ direction of (6) clearly implies $L' \subseteq \{\hat{\mathbf{T}} \mid \mathbf{T} \in L\}$.

So suppose $\mathbf{T}' = (T, \ell^{\mathbf{T}'}) \in L'$, and let $\mathbf{T} = (T, \ell^{\mathbf{T}}) = \pi(\mathbf{T}')$. All we need to show is that the equations (4) and (5) hold with \mathbf{T}' in place of $\hat{\mathbf{T}}$.

As for (4), it follows from the fact that $\ell^{\mathbf{T}'}(\varepsilon) \in A'$.

As for (5), suppose $u \in \text{dom}(\prec_m^T)$. Since $(\ell^{\mathbf{T}'}(u), \ell^{\mathbf{T}'}(u \cdot m)) \in K$, $\ell^{\mathbf{T}'}(u \cdot m) = (\mathbf{V}, \varepsilon)$ for some $\mathbf{V} \in F$. We show two things:

$$v \in C_m^T(u) \text{ implies } v \in V \text{ and } \ell^{\mathbf{T}'}(u \cdot m \cdot v) = (\mathbf{V}, v). \quad (7)$$

$$V \subseteq C_m^T(u). \quad (8)$$

It then easily follows that $\mathbf{V} = \mathbf{C}_m^T(u)$ and (5) holds whenever $v \in C_m^T(u)$.

We show (7) by induction on $v \in C_m^T(u)$. For $v = \varepsilon$, we already know that $\varepsilon \in V$ and $\ell^{\mathbf{T}'}(u \cdot m) = (\mathbf{V}, \varepsilon)$. If $v \neq \varepsilon$, we can write $v = v' \cdot (m-1) \cdot v''$ with $v'' \in \mathbb{P}_{m-2}$. Suppose $u \cdot m \cdot v' \prec_{m-1,i}^T u \cdot m \cdot v$. Since $v' \in C_m^T(u)$, by induction hypothesis, $v' \in V$ and $\ell^{\mathbf{T}'}(u \cdot m \cdot v') = (\mathbf{V}, v')$. Since $\mathbf{T}' \in L'$, $(\ell^{\mathbf{T}'}(u \cdot m \cdot v'), C_{m-1}^T(u \cdot m \cdot v'), i, \ell^{\mathbf{T}'}(u \cdot m \cdot v)) \in J$. By the definition of J , we must have $C_{m-1}^T(u \cdot m \cdot v') = C_{m-1}^V(v')$, which implies $v' \prec_{m-1,i}^V v \in V$. The definition of J then implies $\ell^{\mathbf{T}'}(u \cdot m \cdot v) = (\mathbf{V}, v)$.

Having established (7), we proceed to show (8) by induction on $v \in V$. For $v = \varepsilon$, we have $\varepsilon \in C_m^T(u)$ since $u \in \text{dom}(\prec_m^T)$. If $v = v' \cdot (m-1) \cdot v''$ with $v'' \in \mathbb{P}_{m-2}$, then $v' \in C_m^T(u)$ by induction hypothesis. Since $v' \in \text{dom}(\prec_{m-1}^V)$, $\ell^{\mathbf{T}'}(u \cdot m \cdot v') = (\mathbf{V}, v') \notin Y$. By the assumption that $\mathbf{T}' \in L'$, this means that $u \cdot m \cdot v' \in \text{dom}(\prec_{m-1}^T)$ and so $(\ell^{\mathbf{T}'}(u \cdot m \cdot v'), C_{m-1}^T(u \cdot m \cdot v'), 1, \ell^{\mathbf{T}'}(u \cdot m \cdot v' \cdot (m-1))) \in J$. Since $\ell^{\mathbf{T}'}(u \cdot m \cdot v') = (\mathbf{V}, v')$, the definition of J implies $C_{m-1}^T(u \cdot m \cdot v') = C_{m-1}^V(v')$. Since $v = v' \cdot (m-1) \cdot v'' \in V$, it follows that $v'' \in C_{m-1}^T(u \cdot m \cdot v')$ and hence $v = v' \cdot (m-1) \cdot v'' \in C_m^T(u)$.

This concludes the proof of the lemma. \square

6 Encoding and Yield at Higher Dimensions

In order to prove an analogue of the Chomsky-Schützenberger theorem for the m -dimensional yields of local $(m+1)$ -dimensional tree languages, we need to define the higher-dimensional counterparts of the mappings \mathbf{enc}, \mathbf{y} , and of the Dyck languages. Since we use 3-dimensional trees to represent derivation trees

of simple context-free tree grammars, the yield function mapping 3-dimensional trees to 2-dimensional trees must be consistent with the relation between derivation trees of simple context-free tree grammars and their tree yield.

We set aside a special symbol \mathbf{x} and use it to extend a given set Σ of symbols. The intended role of \mathbf{x} in m -dimensional trees over $\Sigma \cup \{\mathbf{x}\}$ is that of a placeholder; the encoding function erases all occurrences of \mathbf{x} . We write $\mathbb{T}_{\Sigma}^m(n)$ to denote the set of m -dimensional trees in $\mathbb{T}_{\Sigma \cup \{\mathbf{x}\}}^m$ in which \mathbf{x} labels exactly n nodes and none of these nodes have a child in the m -th dimension. Let $\mathbf{T} \in \mathbb{T}_{\Sigma}^m(n)$, $\mathbf{T}_1, \dots, \mathbf{T}_n \in \mathbb{T}_{\Sigma \cup \{\mathbf{x}\}}^m$, and let u_1, \dots, u_n be the nodes of \mathbf{T} labeled by \mathbf{x} , in the alphabetical order. Then we write

$$\mathbf{T}[\mathbf{T}_1, \dots, \mathbf{T}_n]$$

for the tree $\mathbf{T}' \in \mathbb{T}_{\Sigma \cup \{\mathbf{x}\}}^m$ such that

$$\begin{aligned} \mathbf{T}' &= \mathbf{T} \cup u_1 \cdot \mathbf{T}_1 \cup \dots \cup u_n \cdot \mathbf{T}_n, \\ \ell^{\mathbf{T}'}(v) &= \begin{cases} \ell^{\mathbf{T}}(v) & \text{if } v \in \mathbf{T} - \{u_1, \dots, u_n\}, \\ \ell^{\mathbf{T}_i}(v') & \text{if } v = u_i \cdot v'. \end{cases} \end{aligned}$$

Given an m -dimensional hedge $\mathbf{T} \in \mathbb{H}_{\Sigma \cup \{\mathbf{x}\}}^m$, define a binary relation $\triangleleft_{m,i}^{\mathbf{T}}$ on T for each positive integer i , as follows: $u \triangleleft_{m,i}^{\mathbf{T}} v$ if and only if v is alphabetically the i -th node in $\{w \mid u \triangleleft_m w, \ell^{\mathbf{T}}(w) = \mathbf{x}\}$.

Let $m \geq 2$. An m -dimensional hedge $\mathbf{T} \in \mathbb{H}_{\Sigma \cup \{\mathbf{x}\}}^m$ is *well-labeled* if

- for all $v \in T$, $\ell^{\mathbf{T}}(v) = \mathbf{x}$ implies $v \notin \text{dom}(\triangleleft_m^{\mathbf{T}}) \cup \text{dom}(\triangleleft_{m-1}^{\mathbf{T}})$, and
- for all $v \in \text{dom}(\triangleleft_m^{\mathbf{T}})$, $\mathbf{C}_m^{\mathbf{T}}(v) \in \mathbb{T}_{\Sigma}^{m-1}(n)$ implies $|C_{m-1}(v)| = n$.

We write $\mathbb{H}_{\Sigma, \mathbf{x}}^m$ to denote the class of well-labeled m -dimensional hedges over $\Sigma \cup \{\mathbf{x}\}$. If $\mathbf{T} \in \mathbb{H}_{\Sigma, \mathbf{x}}^m$, then for each node $v \in \text{dom}(\triangleleft_m^{\mathbf{T}})$, there is a bijection between $\{u \in T \mid v \triangleleft_m^{\mathbf{T}} u \text{ and } \ell^{\mathbf{T}}(u) = \mathbf{x}\}$ and $\{u \in T \mid v \triangleleft_{m-1}^{\mathbf{T}} u\}$, namely, $\bigcup_{i \geq 1} ((\triangleleft_{m,i}^{\mathbf{T}})^{-1} \circ \triangleleft_{m-1,i}^{\mathbf{T}})$. We write $\mathbb{H}_{\Sigma, \mathbf{x}}^m(n)$ for

$$\{\mathbf{T} \in \mathbb{H}_{\Sigma, \mathbf{x}}^m \mid \text{there are exactly } n \text{ nodes } v \in T \cap \mathbb{P}_{m-1} \text{ such that } \ell^{\mathbf{T}}(v) = \mathbf{x}\}.$$

Note that if $\mathbf{T} \in \mathbb{H}_{\Sigma, \mathbf{x}}^m(n)$, \mathbf{T} may have more than n nodes labeled by \mathbf{x} .

We write $\mathbb{T}_{\Sigma, \mathbf{x}}^m$ for $\mathbb{H}_{\Sigma, \mathbf{x}}^m(0) \cap \mathbb{T}_{\Sigma \cup \{\mathbf{x}\}}^m$. We will give suitable definitions of encoding and yield for elements of $\mathbb{T}_{\Sigma, \mathbf{x}}^m$ shortly. Before that, here is a variant of Lemma 10 for languages consisting of well-labeled m -dimensional trees. If $\pi: \Sigma' \rightarrow \Sigma$ is a projection, we extend it to a projection $\pi: \Sigma' \cup \{\mathbf{x}\} \rightarrow \Sigma \cup \{\mathbf{x}\}$ by letting $\pi(\mathbf{x}) = \mathbf{x}$.

Lemma 11. *Let $m \geq 2$. For any local m -dimensional tree language $L \subseteq \mathbb{T}_{\Sigma, \mathbf{x}}^m$, there exist a finite alphabet Σ' , a degree-bounded, super-local m -dimensional tree language $L' \subseteq \mathbb{T}_{\Sigma', \mathbf{x}}^m$, and a projection $\pi: \Sigma' \rightarrow \Sigma$ such that $L = \pi(L')$. Moreover, π maps L' bijectively L' to L .*

Proof. Since the case where $L \subseteq \mathbb{T}_{\Sigma}^m$ is covered by Lemma 10, we assume $L \not\subseteq \mathbb{T}_{\Sigma}^m$. Without loss of generality, we may assume that $L = \text{Loc}^m(A, Z, I)$, for some $A \subseteq \Sigma$, $Z \subseteq \Sigma \cup \{\mathbf{x}\}$, $I \subseteq \Sigma \times \mathbb{T}_{\Sigma \cup \{\mathbf{x}\}}^{m-1}$ such that

- $\mathbf{x} \in Z$,
- $(c, \mathbf{T}) \in I$ implies $c \in \Sigma$ and $\ell^{\mathbf{T}}(v) \in \Sigma$ for all $v \in \text{dom}(\prec_{m-1}^{\mathbf{T}})$, and
- there exist $(c, \mathbf{T}) \in I$ and $v \in T$ such that $\ell^{\mathbf{T}}(v) = \mathbf{x}$.

We modify the construction in the proof of Lemma 10 slightly. The difference is that where (\mathbf{T}, v) would appear in the earlier construction, \mathbf{x} appears instead just in case $\ell^{\mathbf{T}}(v) = \mathbf{x}$. Otherwise, the proof is essentially the same.

The definition of Σ'' is changed as follows:

$$\Sigma'' = \{ (\mathbf{T}, v) \mid \mathbf{T} \in \mathbb{T}_{\Sigma \cup \{\mathbf{x}\}}^{m-1}, v \in T, \ell^{\mathbf{T}}(v) \in \Sigma \}.$$

The definition of $\pi: \Sigma'' \rightarrow \Sigma$ remains the same:

$$\pi((\mathbf{T}, v)) = \ell^{\mathbf{T}}(v).$$

As before, let

$$\begin{aligned} F &= \{ \mathbf{T}_c \mid c \in A \} \cup \{ \mathbf{T} \mid (c, \mathbf{T}) \in I \}, \\ \Sigma' &= \{ (\mathbf{T}, v) \in \Sigma'' \mid \mathbf{T} \in F \}. \end{aligned}$$

The definitions of Z', K, Y, J are modified as follows:

$$\begin{aligned} A' &= \{ (\mathbf{T}_c, \varepsilon) \mid c \in A \}, \\ Z' &= \{ \mathbf{x} \} \cup \{ (\mathbf{T}, v) \in \Sigma' \mid \ell^{\mathbf{T}}(v) \in Z - \{ \mathbf{x} \} \}, \\ K &= \{ ((\mathbf{T}, v), (\mathbf{T}', \varepsilon)) \mid (\mathbf{T}, v) \in \Sigma', (\ell^{\mathbf{T}}(v), \mathbf{T}') \in I, \ell^{\mathbf{T}'}(\varepsilon) \in \Sigma \} \cup \\ &\quad \{ ((\mathbf{T}, v), \mathbf{x}) \mid (\mathbf{T}, v) \in \Sigma', (\ell^{\mathbf{T}}(v), \mathbf{T}_{\mathbf{x}}) \in I \}, \\ Y &= \{ \mathbf{x} \} \cup \{ (\mathbf{T}, v) \in \Sigma' \mid v \notin \text{dom}(\prec_{m-1}^{\mathbf{T}}) \}, \\ J &= \{ ((\mathbf{T}, u), C_{m-1}^{\mathbf{T}}(u), i, (\mathbf{T}, v)) \mid (\mathbf{T}, u) \in \Sigma', u \triangleleft_{m-1, i}^{\mathbf{T}} v, \ell^{\mathbf{T}}(v) \in \Sigma \} \cup \\ &\quad \{ ((\mathbf{T}, u), C_{m-1}^{\mathbf{T}}(u), i, \mathbf{x}) \mid (\mathbf{T}, u) \in \Sigma', u \triangleleft_{m-1, i}^{\mathbf{T}} v, \ell^{\mathbf{T}}(v) = \mathbf{x} \}. \end{aligned}$$

These are finite sets. As before, let $L' \subseteq \mathbb{T}_{\Sigma'' \cup \{\mathbf{x}\}}^m$ be the super-local set defined by $L' = \text{SLoc}(A', Z', K, Y, J)$. It is easy to see that $L' \subseteq \mathbb{T}_{\Sigma' \cup \{\mathbf{x}\}}^m$.

For each $\mathbf{T} \in \mathbb{T}_{\Sigma \cup \{\mathbf{x}\}}^m$, define an m -dimensional tree $\hat{\mathbf{T}} = (T, \ell^{\hat{\mathbf{T}}}) \in \mathbb{T}_{\Sigma'' \cup \{\mathbf{x}\}}^m$ by

$$\ell^{\hat{\mathbf{T}}}(\varepsilon) = (\mathbf{T}_{\ell^{\mathbf{T}}(\varepsilon)}, \varepsilon), \tag{9}$$

$$\ell^{\hat{\mathbf{T}}}(u \cdot m \cdot v) = \begin{cases} (C_m^{\mathbf{T}}(u), v) & \text{if } \ell^{\mathbf{T}}(u \cdot m \cdot v) \in \Sigma, \\ \mathbf{x} & \text{if } \ell^{\mathbf{T}}(u \cdot m \cdot v) = \mathbf{x}, \end{cases} \quad \text{for } v \in C_{m-1}^{\mathbf{T}}(u). \tag{10}$$

It is clear that $\pi(\hat{\mathbf{T}}) = \mathbf{T}$ for all $\mathbf{T} \in \mathbb{T}_{\Sigma \cup \{\mathbf{x}\}}^m$. Similarly to Lemma 10, we can show

$$L' = \{ \hat{\mathbf{T}} \mid \mathbf{T} \in L \}.$$

This clearly implies that π is a bijection from L' to L . \square

Fix $m \geq 2$ and Σ . For each $c \in \Sigma$ and each (possibly empty) finite prefix-closed subset P of \mathbb{P}_{m-1} , define

$$\Gamma_{c,P} = \{ (c, P, i) \mid 0 \leq i \leq |P| \}.$$

We consider $\Gamma_{c,P}$ to be a group of symbols that match with each other; this notion of a matching group of symbols generalizes the notion of a matching pair of brackets. Let

$$\tilde{\Sigma} = \Sigma \cup \bigcup \{ \Gamma_{c,P} \mid c \in \Sigma \text{ and } P \subseteq \mathbb{P}_{m-1} \text{ is finite and prefix-closed} \}.$$

Note that $\tilde{\Sigma}$ is an infinite set.²⁰

Let $\mathbf{T} \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}$. For $u \in C_m^T(\varepsilon)$, let $T_u = \{ v \in \mathbb{P}_m \cdot \mathbb{P}_{m+1} \mid m \cdot u \cdot v \in T \}$, i.e., the domain of $SH_{m+1}(\mathbf{T}, m \cdot u)$. Then we have

$$T = \begin{cases} \{\varepsilon\} \cup \bigcup_{u \in C_m^T(\varepsilon)} m \cdot u \cdot T_u & \text{if } m+1 \notin T, \\ \{\varepsilon\} \cup (m+1) \cdot (T/(m+1)) \cup \bigcup_{u \in C_m^T(\varepsilon)} m \cdot u \cdot T_u & \text{if } m+1 \in T. \end{cases}$$

Thus, \mathbf{T} is completely determined by the following pieces of information:

- $\ell^{\mathbf{T}}(\varepsilon)$,
- $C_m^T(\varepsilon)$,
- $SH_{m+1}(\mathbf{T}, m \cdot u)$ for each $u \in C_m^T(\varepsilon)$,
- whether or not $m+1 \in T$, and
- in case $m+1 \in T$, the $(m+1)$ -dimensional hedge $SH_{m+1}(\mathbf{T}, m+1) = \mathbf{T}/(m+1)$.

Let $P = C_m^T(\varepsilon)$, $k = |P|$, and for $i = 1, \dots, k$, $\varepsilon \triangleleft_{m,i}^T m \cdot u_i$ and $\mathbf{T}_i = SH_{m+1}(\mathbf{T}, m \cdot u_i)$. In case $m+1 \in T$ or $k \geq 1$ (i.e., $m \in T$), let $c = \ell^{\mathbf{T}}(\varepsilon) \in \Sigma$. The m -dimensional encoding of \mathbf{T} , $\mathbf{enc}_m(\mathbf{T})$ in symbols, is defined as follows:

$$\mathbf{enc}_m(\mathbf{T}) = \begin{cases} \mathbf{T}_{\ell^{\mathbf{T}}(\varepsilon)} & \text{if } m+1 \notin T \text{ and } k = 0, \\ c -_m P(\mathbf{enc}_m(\mathbf{T}_1), \dots, \mathbf{enc}_m(\mathbf{T}_k)) & \text{if } m+1 \notin T \text{ and } k \geq 1, \\ (c, P, 0) -_m (\mathbf{enc}_m(\mathbf{T}_0)) [& \text{if } m+1 \in T \text{ and} \\ \quad (c, P, 1) -_m \mathbf{enc}_m(\mathbf{T}_1), \quad \mathbf{T}_0 = SH_{m+1}(\mathbf{T}, m+1). \\ \quad \dots, \\ \quad (c, P, k) -_m \mathbf{enc}_m(\mathbf{T}_k) \\ \quad] \end{cases}$$

(The substitution notation in the last clause presupposes $\mathbf{enc}_m(\mathbf{T}_0) \in \mathbb{T}_{\tilde{\Sigma}}^m(k)$, and this is indeed the case as shown by the following lemma.)

²⁰ When we define $\tilde{\Sigma}$ from Σ in this way, we assume that $\Sigma \cap \Gamma_{c,P} = \emptyset$ for all $c \in \Sigma$ and all finite and prefix-closed $P \subseteq \mathbb{P}_{m-1}$. Technically, this assumption may not always be satisfied; nevertheless, we always regard the symbols in $\Gamma_{c,P}$ as “new” symbols. If more rigor is desired, it can be achieved by complicating the definition of $\Gamma_{c,P}$.

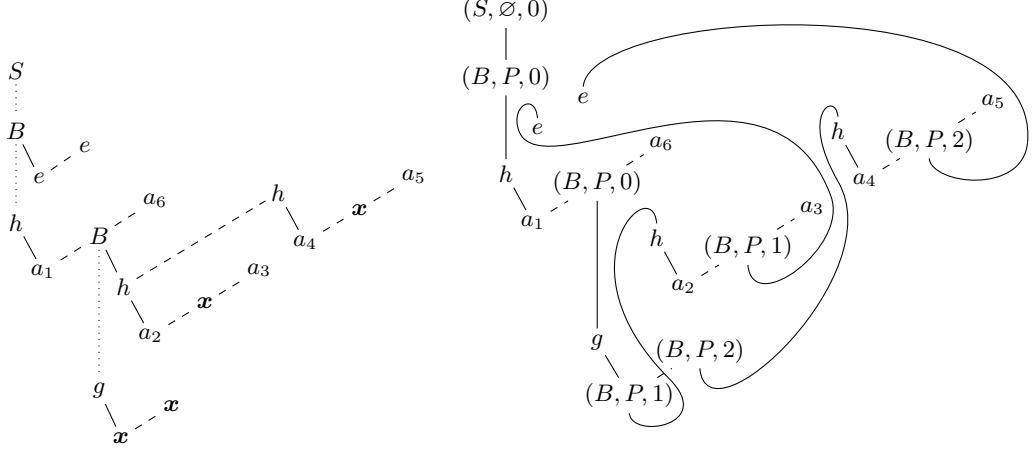


Fig. 3. A well-labeled 3-dimensional tree and its 2-dimensional encoding ($P = \{\varepsilon, 1\}$).

Example 12. Fig. 3 shows the 3-dimensional tree \mathbf{T} from Example 9, which is in $\mathbb{T}_{N \cup \Sigma, \mathbf{x}}^3$, where $N = \{S, B\}$ and $\Sigma = \{h, g, a_1, a_2, a_3, a_4, a_5, a_6\}$, along with $\mathbf{enc}_2(\mathbf{T})$. Here, $P = \{\varepsilon, 1\}$. As before, the edges in the third dimension are dotted, those in the second dimension solid, and those in the first dimension dashed.

Lemma 13. *If $\mathbf{T} \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(n)$, then $\mathbf{enc}_m(\mathbf{T}) \in \mathbb{T}_{\Sigma}^m(n)$.*

Proof. Induction on the size of \mathbf{T} . Let c, P, k , and \mathbf{T}_i be as above. Suppose $\mathbf{T} \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(n)$. If $\ell^{\mathbf{T}}(\varepsilon) = \mathbf{x}$, then $\mathbf{T} = \mathbf{T}_{\mathbf{x}} \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(1)$, and $\mathbf{enc}_m(\mathbf{T}) = \mathbf{T}_{\mathbf{x}} \in \mathbb{T}_{\Sigma}^m(1)$. If $\ell^{\mathbf{T}}(\varepsilon) = c \in \Sigma$, then $\mathbf{T}_i \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(n_i)$ for $i = 1, \dots, k$, where $n = n_1 + \dots + n_k$. By induction hypothesis, $\mathbf{enc}_m(\mathbf{T}_i) \in \mathbb{T}_{\Sigma}^m(n_i)$. Suppose $m+1 \notin T$. Then it is easy to see $\mathbf{enc}_m(\mathbf{T}) \in \mathbb{T}_{\Sigma}^m(n)$. Now suppose $m+1 \in T$. Since \mathbf{T} is well-labeled, $\mathbf{T}_0 \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(k)$. By induction hypothesis, $\mathbf{enc}_m(\mathbf{T}_0) \in \mathbb{T}_{\Sigma}^m(k)$. Also, $(c, P, i) -_m \mathbf{enc}_m(\mathbf{T}_i)$ is in $\mathbb{T}_{\Sigma}^m(n_i)$ for $i = 1, \dots, k$. It easily follows that $\mathbf{enc}_m(\mathbf{T}) = (c, P, 0) -_m (\mathbf{enc}_m(\mathbf{T}_0))[(c, P, 1) -_m \mathbf{enc}_m(\mathbf{T}_1), \dots, (c, P, k) -_m \mathbf{enc}_m(\mathbf{T}_k)] \in \mathbb{T}_{\Sigma}^m(n)$. \square

Note that in $\mathbf{enc}_m(\mathbf{T})$, every node with a label of the form (c, P, i) ($i \geq 0$) has exactly one child in the m -th dimension. There is a simple way of deleting any collection of such nodes from an m -dimensional tree to produce another m -dimensional tree.

Let \mathbf{T} be any m -dimensional tree, and assume that $U \subseteq T$ only consists of nodes v such that $|C_m^T(v)| = 1$. Define a function $f_U: T \rightarrow [1, m]^*$ by

$$\begin{aligned} f_U(\varepsilon) &= \varepsilon, \\ f_U(v \cdot i) &= f_U(v) \cdot i \quad \text{for } i < m, \\ f_U(v \cdot m) &= \begin{cases} f_U(v) \cdot m & \text{if } v \notin U, \\ f_U(v) & \text{if } v \in U. \end{cases} \end{aligned}$$

Let $T' = \text{ran}(f_U) = \{f_U(v) \mid v \in T\}$ and $f'_U = f_U \upharpoonright (T - U)$. Then it is easy to see that $T' = \text{ran}(f'_U)$, T' is a non-empty prefix-closed subset of \mathbb{P}_m , and f'_U is a bijection from $T - U$ to T' . Define

$$\mathbf{del}_m(\mathbf{T}, U) = (T', \ell'),$$

where

$$\ell'(v) = \ell^{\mathbf{T}}((f'_U)^{-1}(v)).$$

Since T' is a non-empty prefix-closed subset of \mathbb{P}_m , it follows that $\mathbf{del}_m(\mathbf{T}, U)$ is an m -dimensional tree.

Let $\mathcal{Y} \subseteq \Sigma$ and $\mathbf{T} \in \mathbb{T}_{\Sigma}^m$. We define

$$\mathbf{del}_{m, \mathcal{Y}}(\mathbf{T}) = \mathbf{del}_m(\mathbf{T}, U),$$

where

$$U = \{v \in T \mid \ell^{\mathbf{T}}(v) \in \mathcal{Y} \text{ and } |C_m^T(v)| = 1\}.$$

Now let $\mathbf{T} \in \mathbb{T}_{\Sigma, \mathbf{x}}^{m+1}$ ($m \geq 2$). The m -dimensional yield of \mathbf{T} is defined as follows:

$$\mathbf{y}_m(\mathbf{T}) = \mathbf{del}_{m, \tilde{\Sigma} - \Sigma}(\mathbf{enc}_m(\mathbf{T})).$$

It is easy to see that $\mathbf{y}_m(\mathbf{T}) \in \mathbb{T}_{\Sigma}^m$.

It is of course straightforward to define $\mathbf{y}_m: \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(n) \rightarrow \mathbb{T}_{\Sigma}^m(n)$ directly:

$$\mathbf{y}_m(\mathbf{T}) = \begin{cases} \mathbf{T}_c & \text{if } m+1 \notin T \text{ and } k = 0, \\ c -_m P(\mathbf{y}_m(\mathbf{T}_1), \dots, \mathbf{y}_m(\mathbf{T}_k)) & \text{if } m+1 \notin T \text{ and } k \geq 1, \\ (\mathbf{y}_m(\mathbf{T}_0))[\mathbf{y}_m(\mathbf{T}_1), \dots, \mathbf{y}_m(\mathbf{T}_k)] & \text{if } m+1 \in T \text{ and} \\ & \mathbf{T}_0 = SH_{m+1}(\mathbf{T}, m+1), \end{cases}$$

where, as before, $c = \ell^{\mathbf{T}}(\varepsilon)$, $P = C_m^T(\varepsilon)$, $k = |P|$, and for $i = 1, \dots, k$, $\varepsilon \triangleleft_{m, i}^T m \cdot u_i$ and $\mathbf{T}_i = SH_{m+1}(\mathbf{T}, m \cdot u_i)$. The indirect definition through \mathbf{enc}_m , however, is useful for our generalization of the Chomsky-Schützenberger theorem for multi-dimensional tree languages.²¹

²¹ Our definition of \mathbf{y}_m is basically a special case of the “frontier” function of Baldwin and Strawn [2]. The difference is that Baldwin and Strawn index variables with node addresses, whereas here we rely on the relative positions of the variables to pair variables with “sub-hedges”.

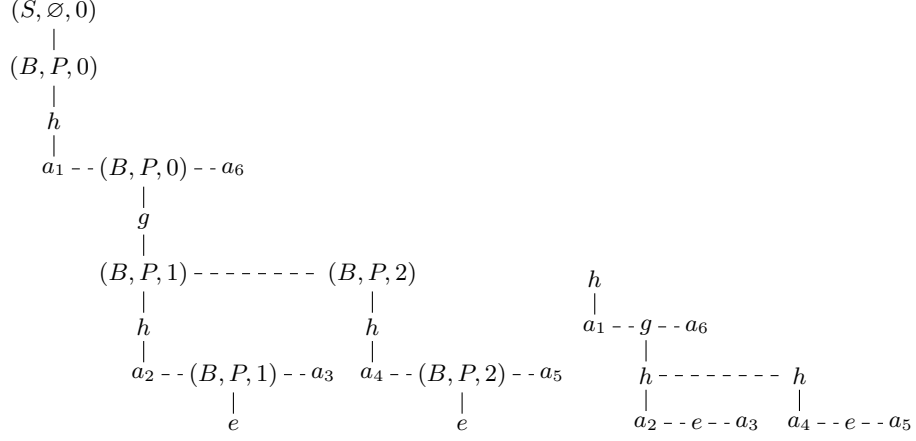


Fig. 4. The 2-dimensional encoding and 2-dimensional yield of a well-labeled 3-dimensional tree ($P = \{\varepsilon, 1\}$).

The case $m = 2$ of the above definition of \mathbf{y}_m is meant to capture the notion of the (tree) yield of a derivation tree of a simple context-free tree grammar, which we represent as a (well-labeled) 3-dimensional tree. The definitions of \mathbf{enc}_m , $\mathbf{del}_{m,\gamma}$, \mathbf{y}_m are all applicable to the case $m = 1$ as well, but the resulting definitions of \mathbf{enc}_1 and of \mathbf{y}_1 will not be equivalent to the standard ones, so we will continue to treat $m = 1$ as a special case.

Example 14. Fig. 4 shows $\mathbf{enc}_2(\mathbf{T})$ (the same tree as the right tree in Fig. 3 with the nodes rearranged) and $\mathbf{y}_2(\mathbf{T})$, where \mathbf{T} is the 3-dimensional tree from Example 9 (the left tree in Fig. 3).

7 Multi-dimensional Dyck Languages

We continue to work with the alphabet

$$\tilde{\Sigma} = \Sigma \cup \bigcup \{ \Gamma_{c,P} \mid c \in \Sigma \text{ and } P \text{ is a finite prefix-closed subset of } \mathbb{P}_{m-1} \},$$

as defined in the previous section. The range of the function $\mathbf{enc}_m : \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(0) \rightarrow \mathbb{T}_{\tilde{\Sigma}}^m$ forms a special subset of $\mathbb{T}_{\tilde{\Sigma}}^m$ similar to Dyck languages.

Let us define a rewriting relation \rightsquigarrow on $\mathbb{T}_{\tilde{\Sigma}}^m$:

$$\mathbf{T} \rightsquigarrow \mathbf{T}'$$

holds if there exist some $v_0, v_1, \dots, v_n \in T$ ($n \geq 0$), $c \in \Sigma$, and finite prefix-closed subset P of \mathbb{P}_m such that²²

²² Note that for $u, v \in T$, $u \langle \triangleleft_m^T \rangle^* v$ is equivalent to $v \in u \cdot \mathbb{P}_m$, and $u \langle \triangleleft_m^T \rangle^+ v$ is equivalent to $v \in u \cdot m \cdot \mathbb{P}_{m-1} \cdot \mathbb{P}_m$.

- $\mathbf{T}' = \mathbf{del}_m(\mathbf{T}, \{v_0, v_1, \dots, v_n\})$,
- $|C_m^T(v_i)| = 1$ for $i = 0, 1, \dots, n$,
- $\ell^{\mathbf{T}}(v_i) = (c, P, i)$ for $i = 0, 1, \dots, n$,
- $n = |P|$,
- v_1, \dots, v_n is the alphabetical listing of $\{v_1, \dots, v_n\}$,
- $v_0 (\triangleleft_m^T)^+ v_i$ for $i = 1, \dots, n$,
- for every $i, j \in \{1, \dots, n\}$, if $v_i (\triangleleft_m^T)^* v_j$, then $i = j$, and
- for every $u \in T$, if $v_0 (\triangleleft_m^T)^+ u$ and there is no $i \in \{1, \dots, n\}$ such that $v_i (\triangleleft_m^T)^* u$, then $\ell^{\mathbf{T}}(u) \in \Sigma$.

Using the term notation, we can write

$\mathbf{T} \rightsquigarrow \mathbf{T}'$ if and only if

$$\mathbf{T} = \mathbf{U}[(c, P, 0) -_m \mathbf{T}_0[(c, P, 1) -_m \mathbf{T}_1, \dots, (c, P, n) -_m \mathbf{T}_n]],$$

$$\mathbf{T}' = \mathbf{U}[\mathbf{T}_0[\mathbf{T}_1, \dots, \mathbf{T}_n]]$$

for some $\mathbf{U} \in \mathbb{T}_\Sigma^m(1)$, $\mathbf{T}_0 \in \mathbb{T}_\Sigma^m(n)$, $\mathbf{T}_i \in \mathbb{T}_\Sigma^m$ ($i = 1, \dots, n$),

$c \in \Sigma$, finite and prefix-closed $P \subseteq \mathbb{P}_{m-1}$ with $|P| = n$.

Define the *m-dimensional Dyck tree language* over Σ by²³

$$DT_\Sigma^m = \{ \mathbf{T} \in \mathbb{T}_\Sigma^m \mid \mathbf{T} \rightsquigarrow^* \mathbf{T}' \in \mathbb{T}_\Sigma^m \}.$$

Note that the alphabet of DT_Σ^m (i.e., the set of labels that appear in elements of DT_Σ^m) is infinite.

Just like the ordinary Dyck language D_n of strings over Γ_n has an alternative inductive definition in terms of a context-free grammar, so too the *m-dimensional tree language* DT_Σ^m admits an inductive definition. First, let us extend the definition of \rightsquigarrow to a rewriting relation on $\mathbb{T}_\Sigma^m(n)$ by taking the exact same definition as before, requiring $\ell^{\mathbf{T}}(u) \in \Sigma$ rather than $\ell^{\mathbf{T}}(u) \in \Sigma \cup \{\mathbf{x}\}$ in the consequent of the last condition. Note that this relation is (strongly) confluent:

Lemma 15. *Let $\mathbf{T}, \mathbf{T}_1, \mathbf{T}_2 \in \mathbb{T}_\Sigma^m(n)$. If $\mathbf{T} \rightsquigarrow \mathbf{T}_1$ and $\mathbf{T} \rightsquigarrow \mathbf{T}_2$, then there exists some $\mathbf{T}' \in \mathbb{T}_\Sigma^m(n)$ such that $\mathbf{T}_1 \rightsquigarrow \mathbf{T}'$ and $\mathbf{T}_2 \rightsquigarrow \mathbf{T}'$.*

For each $n \in \mathbb{N}$, we define $DT_\Sigma^m(n)$ by

$$DT_\Sigma^m(n) = \{ \mathbf{T} \in \mathbb{T}_\Sigma^m(n) \mid \mathbf{T} \rightsquigarrow^* \mathbf{T}' \in \mathbb{T}_\Sigma^m(n) \}.$$

Clearly, $DT_\Sigma^m(0) = DT_\Sigma^m$.

Then we can prove that $(X_n)_{n \in \mathbb{N}} = (DT_\Sigma^m(n))_{n \in \mathbb{N}}$ is the unique solution to a certain equation. For $n \in \mathbb{N}$, let X_n be a variable ranging over the subsets of $\mathbb{T}_\Sigma^m(n)$. Consider the following conditions:

²³ For dimension $m = 2$, analogous notions of Dyck tree language have been proposed by Matsubara and Kasai [28] and by Arnold and Dauchet [1] to capture the tree languages generated by tree-adjointing grammars and by (general) context-free tree grammars, respectively.

- C1. $n = 0$ and $\mathbf{T} = \mathbf{T}_c$ for some $c \in \Sigma$.
 C2. $n = 1$ and $\mathbf{T} = \mathbf{T}_x$.
 C3. There exist $k \geq 1$, $n_1, \dots, n_k \geq 0$, $c \in \Sigma$, some finite prefix-closed $P \subseteq \mathbb{P}_{m-1}$, and some $\mathbf{T}_1 \in X_{n_1}, \dots, \mathbf{T}_k \in X_{n_k}$ such that

$$\begin{aligned} n &= \sum_{i=1}^k n_i, \\ |P| &= k, \\ \mathbf{T} &= c -_m P(\mathbf{T}_1, \dots, \mathbf{T}_k). \end{aligned}$$

- C4. There exist $k \geq 0$, $n_1, \dots, n_k \geq 0$, $c \in \Sigma$, some finite prefix-closed $P \subseteq \mathbb{P}_{m-1}$, and some $\mathbf{T}_1 \in X_{n_1}, \dots, \mathbf{T}_k \in X_{n_k}, \mathbf{T}_0 \in X_k$ such that

$$\begin{aligned} n &= \sum_{i=1}^k n_i, \\ |P| &= k, \\ \mathbf{T} &= (c, P, 0) -_m \mathbf{T}_0[(c, P, 1) -_m \mathbf{T}_1, \dots, (c, P, k) -_m \mathbf{T}_k]. \end{aligned}$$

Note that C1 and C2 are conditions on n and \mathbf{T} , while C3 and C4 are conditions on n , \mathbf{T} , and $(X_i)_{i \in \mathbb{N}}$.

Theorem 16. $(X_n)_{n \in \mathbb{N}} = (DT_{\Sigma}^m(n))_{n \in \mathbb{N}}$ is the unique solution to the following biconditional:

$$\mathbf{T} \in X_n \iff (\text{C1} \vee \text{C2} \vee \text{C3} \vee \text{C4}) \quad (11)$$

Proof. We first show that $(X_n)_{n \in \mathbb{N}} = (DT_{\Sigma}^m(n))_{n \in \mathbb{N}}$ satisfies the biconditional (11).

(\Leftarrow). This direction is easily proved by induction.

(\Rightarrow). Suppose $\mathbf{T} \in DT_{\Sigma}^m(n)$. If $|\mathbf{T}| = 1$, then clearly, either C1 or C2 holds.

If $\ell^{\mathbf{T}}(\varepsilon) \in \Sigma$ and $C_m^{\mathbf{T}}(\varepsilon) = P \neq \emptyset$, let $k = |P|$. Then $\mathbf{T} = c -_m P(\mathbf{T}_1, \dots, \mathbf{T}_k)$ for some $\mathbf{T}_1 \in \mathbb{T}_{\Sigma}^m(n_1), \dots, \mathbf{T}_k \in \mathbb{T}_{\Sigma}^m(n_k)$ such that $\sum_{i=1}^k n_i = n$. Since $\mathbf{T} \rightsquigarrow^* \mathbf{T}'$ for some $\mathbf{T}' \in \mathbb{T}_{\Sigma}^m(n)$, it is clear that for $i = 1, \dots, k$, $\mathbf{T}_i \rightsquigarrow^* \mathbf{T}'_i$ for some $\mathbf{T}'_i \in \mathbb{T}_{\Sigma}^m(n_i)$, and hence $\mathbf{T}_i \in DT_{\Sigma}^m(n_i)$. Therefore, C3 holds.

Now suppose $\ell^{\mathbf{T}}(\varepsilon) = (c, P, i)$. Since $\mathbf{T} \rightsquigarrow^* \mathbf{T}' \in \mathbb{T}_{\Sigma}^m(n)$, it is easy to see that $i = 0$ and $C_m^{\mathbf{T}}(\varepsilon) = \{\varepsilon\}$. Let $k = |P|$ and let $\mathbf{T}'_0 \in \mathbb{T}_{\Sigma}^m(k), \mathbf{T}'_1, \dots, \mathbf{T}'_k$ be such that

$$\begin{aligned} \mathbf{T} &\rightsquigarrow^* (c, P, 0) -_m \mathbf{T}'_0[(c, P, 1) -_m \mathbf{T}'_1, \dots, (c, P, k) -_m \mathbf{T}'_k] \\ &\rightsquigarrow \mathbf{T}'_0[\mathbf{T}'_1, \dots, \mathbf{T}'_k] \\ &\rightsquigarrow^* \mathbf{T}'. \end{aligned}$$

Then it is easy to see that for $i = 1, \dots, k$, $\mathbf{T}'_i \rightsquigarrow^* \mathbf{T}''_i \in \mathbb{T}_{\Sigma}^m(n_i)$ for some n_i such that $n = \sum_{i=1}^k n_i$. Also, we must have

$$\mathbf{T} = (c, P, 0) -_m \mathbf{T}_0[(c, P, 1) -_m \mathbf{T}_1, \dots, (c, P, k) -_m \mathbf{T}_k]$$

for some $\mathbf{T}_0 \in \mathbb{T}_{\Sigma}^m(k)$, $\mathbf{T}_1 \in \mathbb{T}_{\Sigma}^m(n_1)$, \dots , $\mathbf{T}_k \in \mathbb{T}_{\Sigma}^m(n_k)$ such that $\mathbf{T}_0 \rightsquigarrow^* \mathbf{T}'_0$ and for $i = 1, \dots, k$, $\mathbf{T}_i \rightsquigarrow^* \mathbf{T}'_i$. It follows that $\mathbf{T}_0 \in DT_{\Sigma}^m(k)$ and for $i = 1, \dots, k$, $\mathbf{T}_i \in DT_{\Sigma}^m(n_i)$, i.e., C4 holds.

We have shown that $(X_n)_{n \in \mathbb{N}} = (DT_{\Sigma}^m(n))_{n \in \mathbb{N}}$ is a solution to (11). The uniqueness follows from the fact that the existential quantification over $\mathbf{T}_0, \mathbf{T}_1, \dots, \mathbf{T}_k$ in C3 and C4 can be restricted to m -dimensional trees with fewer nodes than \mathbf{T} , so that the biconditional (11) amounts to a simultaneous definition of the characteristic functions of X_0, X_1, \dots by induction on the size of \mathbf{T} . \square

Just as in the case of ordinary Dyck languages, the inductive definition of $DT_{\Sigma}^m(n)$ given in Theorem 16 is *unambiguous* in the sense that every $\mathbf{T} \in DT_{\Sigma}^m(n)$ can be written in the form of one of the equations in C1–C4, in exactly one way. This follows from the next lemma:

Lemma 17. *Let $\mathbf{U} \in DT_{\Sigma}^m(k)$, $\mathbf{U}' \in DT_{\Sigma}^m(l)$. If $\mathbf{U}[(c, P, i_1) \text{ } -_m \mathbf{T}_1, \dots, (c, P, i_k) \text{ } -_m \mathbf{T}_k] = \mathbf{U}'[(c, P, j_1) \text{ } -_m \mathbf{T}'_1, \dots, (c, P, j_l) \text{ } -_m \mathbf{T}'_l]$ with $i_1, \dots, i_k, j_1, \dots, j_l \geq 1$, then $\mathbf{U} = \mathbf{U}'$.*

Proof. This can be proved by straightforward induction on the size of \mathbf{U} , using Theorem 16. \square

Lemma 18. $\{\mathbf{enc}_m(\mathbf{T}) \mid \mathbf{T} \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(n)\} = DT_{\Sigma}^m(n)$.

Proof. By Theorem 16, it suffices to show that $(X_n)_{n \in \mathbb{N}} = (\{\mathbf{enc}_m(\mathbf{T}) \mid \mathbf{T} \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(n)\})_{n \in \mathbb{N}}$ satisfies the biconditional (11).

(\Rightarrow). Suppose $\mathbf{T} \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(n)$. We show that n and $\mathbf{enc}_m(\mathbf{T})$ satisfy one of C1–C4 with respect to $(X_n)_{n \in \mathbb{N}} = (\{\mathbf{enc}_m(\mathbf{T}) \mid \mathbf{T} \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(n)\})_{n \in \mathbb{N}}$.

Let $P = C_m^T(\varepsilon)$, and $k = |P|$. Assume that for $i = 1, \dots, k$, $\varepsilon \triangleleft_{m,i}^T m \cdot u_i$ and $\mathbf{T}_i = SH_{m+1}(\mathbf{T}, m \cdot u_i)$. Clearly, for each $i = 1, \dots, k$, $\mathbf{T}_i \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(n_i)$ for some n_i such that $n = \sum_{i=1}^k n_i$.

Case 1. $m+1 \in T$. Then $\ell^{\mathbf{T}}(\varepsilon) = c$ for some $c \in \Sigma$. Let $\mathbf{T}_0 = SH_{m+1}(\mathbf{T}, m+1)$. Then $\mathbf{T}_0 \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(k)$ and we have $\mathbf{enc}_m(\mathbf{T}) = (c, P, 0) \text{ } -_m (\mathbf{enc}_m(\mathbf{T}_0))[(c, P, 1) \text{ } -_m \mathbf{enc}_m(\mathbf{T}_1), \dots, (c, P, k) \text{ } -_m \mathbf{enc}_m(\mathbf{T}_k)] \in DT_{\Sigma}^m(n)$. So n and $\mathbf{enc}_m(\mathbf{T})$ satisfy C4.

Case 2. $m+1 \notin T$. If $P \neq \emptyset$, then $\ell^{\mathbf{T}}(\varepsilon) = c$ for some $c \in \Sigma$, and $\mathbf{enc}_m(\mathbf{T}) = c \text{ } -_m P(\mathbf{enc}_m(\mathbf{T}_1), \dots, \mathbf{enc}_m(\mathbf{T}_k))$, so n and $\mathbf{enc}_m(\mathbf{T})$ satisfy C3. If $P = \emptyset$, then n equals 0 or 1 depending on whether $\ell^{\mathbf{T}}(\varepsilon) = c \in \Sigma$ or $\ell^{\mathbf{T}}(\varepsilon) = \mathbf{x}$. In the former case, $n = 0$ and $\mathbf{enc}_m(\mathbf{T}) = \mathbf{enc}(\mathbf{T}_c) = \mathbf{T}_c$ satisfy C1. In the latter case, $n = 1$ and $\mathbf{enc}_m(\mathbf{T}) = \mathbf{enc}(\mathbf{T}_x) = \mathbf{T}_x$ satisfy C2.

(\Leftarrow). Suppose that n and \mathbf{T} satisfy one C1–C4 with respect to $(X_n)_{n \in \mathbb{N}} = (\{\mathbf{enc}_m(\mathbf{T}) \mid \mathbf{T} \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(n)\})_{n \in \mathbb{N}}$.

If C1 holds, then $n = 0$ and $\mathbf{T} = \mathbf{T}_c = \mathbf{enc}_m(\mathbf{T}_c)$, so $\mathbf{T} \in \{\mathbf{enc}_m(\mathbf{T}') \mid \mathbf{T}' \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(0)\}$.

If C2 holds, then $n = 1$ and $\mathbf{T} = \mathbf{T}_x = \mathbf{enc}_m(\mathbf{T}_x)$, so $\mathbf{T} \in \{\mathbf{enc}_m(\mathbf{T}') \mid \mathbf{T}' \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(1)\}$.

Suppose C3 holds, so that $\mathbf{T} = c -_m P(\mathbf{enc}_m(\mathbf{T}'_1), \dots, \mathbf{enc}_m(\mathbf{T}'_k))$ with $|P| = k$, $\mathbf{T}'_i \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(n_i)$, $n = \sum_{i=1}^k n_i$. Let u_1, \dots, u_k be the elements of P , in alphabetical order. Define \mathbf{T}' by

$$\begin{aligned} T' &= \{\varepsilon\} \cup \bigcup_{i=1}^k m \cdot u_i \cdot T'_i, \\ \ell^{\mathbf{T}'}(\varepsilon) &= c, \\ \ell^{\mathbf{T}'}(m \cdot u_i \cdot v) &= \ell^{\mathbf{T}'_i}(v) \quad \text{for } v \in T'_i. \end{aligned}$$

Then $\mathbf{T}' \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(n)$ and $\mathbf{T}'_i = SH_{m+1}(\mathbf{T}', m \cdot u_i)$ for $i = 1, \dots, k$. By the definition of \mathbf{enc}_m ,

$$\mathbf{enc}_m(\mathbf{T}') = c -_m P(\mathbf{enc}_m(\mathbf{T}'_1), \dots, \mathbf{enc}_m(\mathbf{T}'_k)) = \mathbf{T}.$$

So $\mathbf{T} \in \{\mathbf{enc}_m(\mathbf{T}') \mid \mathbf{T}' \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(n)\}$.

Suppose C4 holds, so that $\mathbf{T} = (c, P, 0) -_m \mathbf{T}'_0[(c, P, 1) -_m \mathbf{T}'_1, \dots, (c, P, k) -_m \mathbf{T}'_k]$ with $|P| = k$, $\mathbf{T}'_0 \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(k)$, and $\mathbf{T}'_i \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(n_i)$, and $n = \sum_{i=1}^k n_i$. Let u_1, \dots, u_k list the elements of P , in alphabetical order. Define \mathbf{T}' by

$$\begin{aligned} T' &= \{\varepsilon\} \cup (m+1) \cdot T'_0 \cup \bigcup_{i=1}^k m \cdot u_i \cdot T'_i, \\ \ell^{\mathbf{T}'}(\varepsilon) &= c, \\ \ell^{\mathbf{T}'}((m+1) \cdot v) &= \ell^{\mathbf{T}'_0}(v) \quad \text{for } v \in T'_0, \\ \ell^{\mathbf{T}'}(m \cdot u_i \cdot v) &= \ell^{\mathbf{T}'_i}(v) \quad \text{for } v \in T'_i. \end{aligned}$$

Then it is easy to see $\mathbf{T}' \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(n)$, $\mathbf{T}'_0 = SH_{m+1}(\mathbf{T}', m+1)$, and for $i = 1, \dots, k$, $\mathbf{T}'_i = SH_{m+1}(\mathbf{T}', m \cdot u_i)$. By the definition of \mathbf{enc}_m ,

$$\begin{aligned} \mathbf{enc}_m(\mathbf{T}') &= \\ &= (c, P, 0) -_m (\mathbf{enc}_m(\mathbf{T}'_0))[(c, P, 1) -_m (\mathbf{enc}_m(\mathbf{T}'_1)), \dots, (c, P, k) -_m (\mathbf{enc}_m(\mathbf{T}'_k))] \\ &= \mathbf{T}. \end{aligned}$$

So $\mathbf{T} \in \{\mathbf{enc}_m(\mathbf{T}') \mid \mathbf{T}' \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(n)\}$. \square

Lemma 19. For each $m \geq 2$, \mathbf{enc}_m is an injection.

Proof. This follows from the unambiguity of the inductive definition of $DT_{\Sigma}^m(n)$. \square

It is useful to define a function $f_{\mathbf{enc}_m}^{\mathbf{T}}$ from the nodes of $\mathbf{T} \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}$ to the nodes of $\mathbf{T}' = \mathbf{enc}_m(\mathbf{T})$. Let $P = C_m^{\mathbf{T}}(\varepsilon)$ and $k = |P|$. Let u_1, \dots, u_k list the elements of P in alphabetical order, and let $\mathbf{T}'_i = SH_{m+1}(\mathbf{T}, m \cdot u_i)$ for $i = 1, \dots, k$. Define $f_{\mathbf{enc}_m}^{\mathbf{T}} : T \rightarrow T'$ by

- (i) $f_{\mathbf{enc}_m}^{\mathbf{T}}(\varepsilon) = \varepsilon$.
 (ii) If $m+1 \in T$ and $\mathbf{T}_0 = SH_{m+1}(\mathbf{T}, m+1)$, then

$$\begin{aligned} f_{\mathbf{enc}_m}^{\mathbf{T}}((m+1) \cdot w) &= m \cdot f_{\mathbf{enc}_m}^{\mathbf{T}_0}(w) && \text{where } w \in T_0, \\ f_{\mathbf{enc}_m}^{\mathbf{T}}(m \cdot u_i \cdot w) &= m \cdot f_{\mathbf{enc}_m}^{\mathbf{T}_0}(v_i) \cdot m \cdot f_{\mathbf{enc}_m}^{\mathbf{T}_i}(w) && \text{where } w \in T_i \text{ and } \varepsilon \triangleleft_{m+1,i}^{\mathbf{T}} v_i. \end{aligned}$$

- (iii) If $m+1 \notin T$, then

$$f_{\mathbf{enc}_m}^{\mathbf{T}}(m \cdot u_i \cdot w) = m \cdot u_i \cdot f_{\mathbf{enc}_m}^{\mathbf{T}_i}(w) \quad \text{where } w \in T_i.$$

It is easy to check that $f_{\mathbf{enc}_m}^{\mathbf{T}}(v) \in T'$ indeed holds for all $v \in T$.

Example 20. Consider the 3-dimensional tree \mathbf{T} and its 2-dimensional encoding $\mathbf{enc}_2(\mathbf{T})$, depicted in Fig. 3. In these diagrams, the nodes that are related by $f_{\mathbf{enc}_m}^{\mathbf{T}}$ are placed in roughly the same geometrical positions.

Lemma 21. *Let $\mathbf{T} \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}$ and $\mathbf{T}' = \mathbf{enc}_m(\mathbf{T})$. For each $v \in T$, we have*

$$\ell^{\mathbf{T}'}(f_{\mathbf{enc}_m}^{\mathbf{T}}(v)) = \begin{cases} c & \text{if } v \notin \text{dom}(\prec_{m+1}^{\mathbf{T}}) \text{ and } \ell^{\mathbf{T}}(v) = c \in \Sigma, \\ (c, U, 0) & \text{if } v \in \text{dom}(\prec_{m+1}^{\mathbf{T}}), \ell^{\mathbf{T}}(v) = c, \text{ and } C_m^{\mathbf{T}}(v) = U, \\ (c, U, i) & \text{if } \ell^{\mathbf{T}}(v) = \mathbf{x}, u \triangleleft_{m+1,i}^{\mathbf{T}} v, \ell^{\mathbf{T}}(u) = c, \text{ and } C_m^{\mathbf{T}}(u) = U, \\ \mathbf{x} & \text{if } \ell^{\mathbf{T}}(v) = \mathbf{x} \text{ and } v \in T \cap \mathbb{P}_m. \end{cases}$$

Proof. This is easy to see from the definition of \mathbf{enc}_m . \square

The function $f_{\mathbf{enc}_m}^{\mathbf{T}}$ allows an alternative definition by recursion with respect to the alphabetical order on the nodes of \mathbf{T} .²⁴

Lemma 22. *The function $f_{\mathbf{enc}_m}^{\mathbf{T}}$ satisfies the following equations:*

$$\begin{aligned} f_{\mathbf{enc}_m}^{\mathbf{T}}(\varepsilon) &= \varepsilon, \\ f_{\mathbf{enc}_m}^{\mathbf{T}}(u \cdot (m+1)) &= f_{\mathbf{enc}_m}^{\mathbf{T}}(u) \cdot m, \end{aligned}$$

and for $v \in \mathbb{P}_{m-1}$,

$$f_{\mathbf{enc}_m}^{\mathbf{T}}(u \cdot m \cdot v) = \begin{cases} f_{\mathbf{enc}_m}^{\mathbf{T}}(u) \cdot m \cdot v & \text{if } u \notin \text{dom}(\prec_{m+1}^{\mathbf{T}}), \\ f_{\mathbf{enc}_m}^{\mathbf{T}}(u \cdot (m+1) \cdot w) \cdot m & \text{if } u \in \text{dom}(\prec_{m+1}^{\mathbf{T}}), \\ & u \triangleleft_{m,i}^{\mathbf{T}} v, \text{ and } u \triangleleft_{m+1,i}^{\mathbf{T}} w. \end{cases}$$

Proof. The first equation is true by definition. The remaining two equations can be proved by induction on the length of u . \square

Lemma 23. *For all $\mathbf{T} \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}$, $f_{\mathbf{enc}_m}^{\mathbf{T}}$ is a bijection from the nodes of \mathbf{T} to the nodes of $\mathbf{enc}_m(\mathbf{T})$.*

²⁴ This lemma implies that \mathbf{enc}_m , and consequently \mathbf{y}_m , are MSO-definable transductions mapping $(m+1)$ -ary trees to m -ary trees.

Proof. That $f_{\mathbf{enc}_m}^{\mathbf{T}}$ is injective can be shown by induction with respect to the alphabetical order on T using Lemma 22. Since \mathbf{T} and \mathbf{T}' have the same number of nodes, $f_{\mathbf{enc}_m}^{\mathbf{T}}$ must be a bijection. \square

The m -dimensional counterpart DT'_{Σ}^m of the set of Dyck primes ($m \geq 2$) is defined by

$$DT'_{\Sigma}^m = \{(c, \emptyset, 0) -_m \mathbf{T} \mid c \in \Sigma, \mathbf{T} \in DT_{\Sigma}^m\} \cup \{\mathbf{T}_c \mid c \in \Sigma\}.$$

Lemma 24. *For all $\mathbf{T} \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}(0)$, $\mathbf{T} \in \mathbb{T}_{\Sigma, \mathbf{x}}^{m+1}$ if and only if $\mathbf{enc}_m(\mathbf{T}) \in DT'_{\Sigma}^m$.*

Define a function $\rho: \tilde{\Sigma} \rightarrow \Sigma \cup \{\mathbf{x}\}$ by

$$\begin{aligned} \rho(c) &= c \quad \text{for } c \in \Sigma, \\ \rho((c, U, 0)) &= c, \\ \rho((c, U, i)) &= \mathbf{x} \quad \text{for } 1 \leq i \leq |U|. \end{aligned}$$

Then for every $\mathbf{T} \in \mathbb{H}_{\Sigma, \mathbf{x}}^{m+1}$ and $v \in T$, if $\mathbf{T}' = \mathbf{enc}_m(\mathbf{T})$, we have

$$\rho(\ell^{\mathbf{T}'}(f_{\mathbf{enc}_m}^{\mathbf{T}}(v))) = \ell^{\mathbf{T}}(v).$$

The following is a generalization of Lemma 2 to higher dimensions:

Lemma 25. *Let $L \subseteq \mathbb{T}_{\Sigma, \mathbf{x}}^{m+1}$. If L is super-local, then there exists a local set $L' \subseteq \mathbb{T}_{\Sigma}^m$ such that $\mathbf{enc}_m(L) = L' \cap DT'_{\Sigma}^m$.*

Proof. Let L be a super-local subset of $\mathbb{T}_{\Sigma, \mathbf{x}}^{m+1}$. Without loss of generality, we may suppose that $L = \text{SLoc}^{m+1}(A, Z, K, Y, J)$ for some finite sets

$$\begin{aligned} A &\subseteq \Sigma, \\ Z &\subseteq \Sigma \cup \{\mathbf{x}\}, \\ K &\subseteq \Sigma \times (\Sigma \cup \{\mathbf{x}\}), \\ Y &\subseteq \Sigma \cup \{\mathbf{x}\}, \\ J &\subseteq \Sigma \times \{P \subseteq \mathbb{P}_{m-1} \mid P \text{ is an } (m-1)\text{-ary tree domain}\} \times \mathbb{N}_+ \times (\Sigma \cup \{\mathbf{x}\}). \end{aligned}$$

Let

$$\begin{aligned} \Sigma' &= (Z \cap \Sigma) \cup \bigcup \{ \Gamma_{c,U} \mid \mathbf{x} \in Y \cap Z, (c, U, i, a) \in J, (c, b) \in K \} \cup \\ &\quad \{(c, \emptyset, 0) \mid c \in A \cup Y, (c, a) \in K \}. \end{aligned}$$

Note that Σ' is a finite subset of $\tilde{\Sigma}$. Define finite sets $A', Z' \subseteq \Sigma'$ and $I \subseteq \Sigma' \times \mathbb{T}_{\Sigma'}^{m-1}$ by

$$\begin{aligned} A' &= (A \cap Z) \cup \{ (c, \emptyset, 0) \mid c \in A, (c, a) \in K \}, \\ Z' &= (A \cup Y) \cap Z \cap \Sigma, \\ I &= \{ ((c, U, 0), \mathbf{T}_d) \mid (c, U, 0) \in \Sigma', d \in Z \cap \Sigma, (c, d) \in K \} \cup \\ &\quad \{ ((c, U, 0), \mathbf{T}_{(d, V, 0)}) \mid (c, U, 0), (d, V, 0) \in \Sigma', (c, d) \in K, \text{ and} \\ &\quad \quad \text{either } V \neq \emptyset \text{ or } d \in Y \} \cup \\ &\quad \{ ((c, U, 0), \mathbf{T}_{(c, U, 1)}) \mid (c, U, 1) \in \Sigma', |U| = 1, (c, \mathbf{x}) \in K \} \cup \\ &\quad \{ ((c, U, i), \mathbf{T}_d) \mid (c, U, i) \in \Sigma', (c, U, i, d) \in J, d \in Z \cap \Sigma \} \cup \\ &\quad \{ ((c, U, i), \mathbf{T}_{(d, V, 0)}) \mid (c, U, i) \in \Sigma', (c, U, i, d) \in J, (d, V, 0) \in \Sigma', \text{ and} \\ &\quad \quad \text{either } V \neq \emptyset \text{ or } d \in Y \} \cup \\ &\quad \{ ((c, U, i), \mathbf{T}_{(d, V, j)}) \mid (c, U, i) \in \Sigma', (c, U, i, \mathbf{x}) \in J, j \geq 1, (d, V, j) \in \Sigma' \} \cup \\ &\quad \{ (c, \mathbf{U}) \mid c \in Z \cap \Sigma, \mathbf{U} \in \mathbb{T}_{\Sigma'}^{m-1}, \\ &\quad \quad \text{for } i = 1, \dots, |U|, \text{ if } u_i \text{ is the } i\text{-th node of } \mathbf{U}, \text{ then} \\ &\quad \quad (c, U, i, \rho(\ell^{\mathbf{U}}(u_i))) \in J, \text{ and } \ell^{\mathbf{U}}(u_i) = (d, \emptyset, 0) \text{ implies } d \in Y \}. \end{aligned}$$

It is tedious, but not difficult to show that $L' = \text{Loc}^m(A', Z', I)$ satisfies the desired property. We omit the details. \square

The converse of the above lemma does not hold for a reason similar to the one for the case of the standard **enc** function for dimension 1.²⁵

A projection $\pi: \Sigma' \rightarrow \Sigma$ naturally induces a projection $\tilde{\pi}: \tilde{\Sigma}' \rightarrow \tilde{\Sigma}$ in an obvious way:

$$\begin{aligned} \tilde{\pi}(c) &= \pi(c), \\ \tilde{\pi}((c, P, i)) &= (\pi(c), P, i). \end{aligned}$$

Clearly, if $\mathbf{T}' \in DT_{\Sigma'}^m$, then $\tilde{\pi}(\mathbf{T}') \in DT_{\Sigma}^m$. Also, if $\mathbf{T}' \in \mathbb{T}_{\Sigma', \mathbf{x}}^{m+1}$, then $\mathbf{enc}_m(\pi(\mathbf{T}')) = \tilde{\pi}(\mathbf{enc}_m(\mathbf{T}'))$.

Here is a generalization of Lemma 6 to multi-dimensional Dyck languages:

Lemma 26. *Let $L \subseteq \mathbb{T}_{\Sigma}^m$ be a local set. Then there exist a finite alphabet Σ' , a projection $\pi: \Sigma' \rightarrow \Sigma$, and a local set $L' \subseteq \mathbb{T}_{\Sigma'}^m$ such that $L \cap DT_{\Sigma}^m = \tilde{\pi}(L' \cap DT_{\Sigma'}^m)$. Moreover, $\tilde{\pi}$ maps $L' \cap DT_{\Sigma'}^m$ bijectively to $L \cap DT_{\Sigma}^m$.*

Proof. Let $A, Z \subseteq \tilde{\Sigma}$ and $I \subseteq \tilde{\Sigma} \times \mathbb{T}_{\tilde{\Sigma}}^{m-1}$ be finite sets such that $L = \text{Loc}^m(A, Z, I)$. Since we are interested in the intersection of L and DT_{Σ}^m , we may assume without loss of generality $Z \subseteq \Sigma$. Define

$$\begin{aligned} \Sigma_0 &= Z \cup \{ c \in \Sigma \mid (c, \mathbf{T}) \in I \}, \\ \Sigma' &= \Sigma_0 \cup \{ \bar{c} \mid c \in A \cap Z \}, \end{aligned}$$

²⁵ There is also an additional reason. $L = \{a -_3 a\}$ is not super-local even though $\mathbf{enc}_2(L) = \{(a, \emptyset, 0) -_2 a\}$ is local.

$$\begin{aligned} A' &= \{\bar{c} \mid c \in A \cap Z\} \cup \{(c, \emptyset, 0) \mid c \in \Sigma_0, (c, \emptyset, 0) \in A\}, \\ Z' &= Z \cup \{\bar{c} \mid c \in A \cap Z\}. \end{aligned}$$

Then A' and Z' are finite subsets of $\widetilde{\Sigma}'$. Let $\pi: \Sigma' \rightarrow \Sigma$ be the projection defined by

$$\pi(c) = c, \quad \pi(\bar{d}) = d$$

for each $c \in \Sigma_0$ and $d \in A \cap Z$. Let

$$L' = \text{Loc}^m(A', Z', I).$$

It is easy to see that $L \cap DT_{\Sigma}^m = \widetilde{\pi}(L' \cap DT_{\Sigma'}^m)$ and for each $\mathbf{T} \in L \cap DT_{\Sigma}^m$, there is a unique $\mathbf{T}' \in L'$ such that $\pi(\mathbf{T}') = \mathbf{T}$. \square

Lemma 27. *If $L \subseteq \mathbb{T}_{\Sigma, \mathbf{x}}^{m+1}$ is a local set, then there exist a finite alphabet Σ' , a projection $\pi: \Sigma' \rightarrow \Sigma$, and a local set $L' \subseteq \mathbb{T}_{\Sigma'}^m$ such that*

$$\mathbf{enc}_m(L) = \widetilde{\pi}(L' \cap DT_{\Sigma'}^m).$$

Moreover, $\mathbf{enc}_m^{-1} \circ \widetilde{\pi}$ maps $L' \cap DT_{\Sigma'}^m$ bijectively to L .

Proof. By Lemma 11, there exist a projection $\pi_1: \Sigma_1 \rightarrow \Sigma$ and a super-local $L_1 \subseteq \mathbb{T}_{\Sigma_1, \mathbf{x}}^{m+1}$ such that $L = \pi_1(L_1)$. By Lemma 25, there exist a local set $L_2 \subseteq \mathbb{T}_{\Sigma_1}^m$ such that $\mathbf{enc}_m(L_1) = L_2 \cap DT_{\Sigma_1}^m$. By Lemma 26, there exist a projection $\pi_2: \Sigma' \rightarrow \Sigma_1$ and a local set $L' \subseteq \mathbb{T}_{\Sigma'}^m$ such that $L_2 \cap DT_{\Sigma_1}^m = \widetilde{\pi}_2(L' \cap DT_{\Sigma'}^m)$. So

$$\begin{aligned} \mathbf{enc}_m(L) &= \mathbf{enc}_m(\pi_1(L_1)) \\ &= \widetilde{\pi}_1(\mathbf{enc}_m(L_1)) \\ &= \widetilde{\pi}_1(L_2 \cap DT_{\Sigma_1}^m) \\ &= \widetilde{\pi}_1(\widetilde{\pi}_2(L' \cap DT_{\Sigma'}^m)) \\ &= \widetilde{\pi}(L' \cap DT_{\Sigma'}^m), \end{aligned}$$

where $\pi = \pi_1 \circ \pi_2$. Since π_1 is a bijection from L_1 to L and \mathbf{enc}_m is injective, $\widetilde{\pi}_1$ maps $\mathbf{enc}_m(L_1)$ bijectively to $\mathbf{enc}_m(L)$. Since $\widetilde{\pi}_2$ maps $L' \cap DT_{\Sigma'}^m$ bijectively to $L_2 \cap DT_{\Sigma_1}^m = \mathbf{enc}_m(L_1)$, $\widetilde{\pi} = \widetilde{\pi}_1 \circ \widetilde{\pi}_2$ maps $L' \cap DT_{\Sigma'}^m$ bijectively to $\mathbf{enc}_m(L)$. \square

8 A Multi-dimensional Generalization of the Chomsky-Schützenberger Theorem

Let $m \geq 2$. We call $L \subseteq \mathbb{T}_{\Sigma}^m$ *simple context-free* if there exist a finite alphabet Υ and a local set $K \subseteq \mathbb{T}_{\Upsilon, \mathbf{x}}^{m+1}$ such that $L = \mathbf{y}_m(K)$.

For a finite alphabet Σ and $r \geq 0$, we define the finite alphabet

$$\widetilde{\Sigma}_r = \Sigma \cup \bigcup \{ \Gamma_{c,P} \mid c \in \Sigma, P \text{ is finite and prefix-closed, } |P| \leq r \}.$$

For any alphabet \mathcal{Y} and $p \geq 1$, let

$$\mathbb{T}_{\mathcal{Y},p}^m = \{ \mathbf{T} \in \mathbb{T}_{\mathcal{Y}}^m \mid |C_m^T(v)| \leq p \text{ for all } v \in T \}.$$

Clearly, if \mathcal{Y} is finite, $\mathbb{T}_{\mathcal{Y},p}^m$ is a local subset of $\mathbb{T}_{\mathcal{Y}}^m$. Also, any local subset L of $\mathbb{T}_{\mathcal{Y}}^m$ is included in $\mathbb{T}_{\mathcal{Y},p}^m$ for some p , which is just another way of saying L is degree-bounded.

Lemma 28. *Let Σ be a finite set. For $m \geq 2$, $DT_{\Sigma}^m \cap \mathbb{T}_{\Sigma,r,p}^m$ is simple context-free.*

Proof. We adapt the inductive definition of $DT_{\Sigma}^m(n)$ to obtain the required local set. Let $\mathcal{Y} = \tilde{\Sigma}_r \cup \{X_0, \dots, X_r\}$. We write U_k for the set $\{\varepsilon, (m-1), \dots, (m-1)^{k-1}\} \subseteq \mathbb{P}_{m-1}$. Let

$$A_n = \{X_n\} \quad \text{for } n = 0, \dots, r,$$

$$Z = \tilde{\Sigma}_r \cup \{\mathbf{x}\},$$

$$I = \{(X_0, \mathbf{T}_c), (X_1, \mathbf{T}_{\mathbf{x}})\} \cup$$

$$\left\{ \left(\begin{array}{c} c -_m P(\\ X_n, \quad \dots, \\ X_{n_k} -_m U_{n_k}(\mathbf{x}, \dots, \mathbf{x}) \end{array} \right) \left| \begin{array}{l} P \subseteq \mathbb{P}_{m-1}, \\ P \text{ is finite and prefix-closed,} \\ 1 \leq |P| = k \leq p, \\ 0 \leq n = n_1 + \dots + n_k \leq r \end{array} \right. \right\} \cup$$

$$\left\{ \left(\begin{array}{c} (c, P, 0) -_m X_k -_m U_k(\\ X_n, \quad \dots, \\ (c, P, k) -_m X_{n_k} -_m U_{n_k}(\mathbf{x}, \dots, \mathbf{x}) \end{array} \right) \left| \begin{array}{l} P \subseteq \mathbb{P}_{m-1}, \\ P \text{ is finite and prefix-closed,} \\ 0 \leq |P| = k \leq r, \\ 0 \leq n = n_1 + \dots + n_k \leq r \end{array} \right. \right\}.$$

Here, the number of occurrences of \mathbf{x} in $U_{n_i}(\mathbf{x}, \dots, \mathbf{x})$ is $|U_{n_i}| = n_i$. When $j = 0$, we understand the notation $X_j -_m U_j(\mathbf{x}, \dots, \mathbf{x})$ to mean X_0 , i.e., the tree consisting of a single node labeled by X_0 . Note that A_n and Z are (finite) subsets of \mathcal{Y} and I is a finite subset of $\mathcal{Y} \times \mathbb{T}_{\mathcal{Y} \cup \{\mathbf{x}\}}^m$. It is straightforward to prove that

$$DT_{\Sigma}^m(n) \cap \mathbb{T}_{\Sigma,r,p}^m(n) = \mathbf{y}_m(\text{Loc}^{m+1}(A_n, Z, I))$$

holds for $n = 0, \dots, r$. The case of $n = 0$ gives the statement of the lemma. We omit the details. \square

I state the next lemma without proof. Part (ii) and (iii) are straightforward. Part (i) can be proved by using the notion of a recognizable (equivalently, MSO-definable) set of m -dimensional trees [33,32] and relying on the fact that the yield mapping is an MSO-definable transduction.

Lemma 29. *Let $L \subseteq \mathbb{T}_{\Sigma}^m$ be a simple context-free set.*

(i) *If $L' \subseteq \mathbb{T}_{\Sigma}^m$ is local, then $L \cap L'$ is simple context-free.*

- (ii) For every projection $\pi: \Sigma \rightarrow \Sigma'$, $\pi(L)$ is simple context-free.
- (iii) If $\Sigma' \subseteq \Sigma$, then $\mathbf{del}_{m, \Sigma'}(L)$ is simple context-free.

Clearly, $\mathbf{T} \in DT_{\Sigma}^m$ implies $\mathbf{del}_{m, \tilde{\Sigma}-\Sigma}(\mathbf{T}) \in \mathbb{T}_{\Sigma}^m$. We obtain the following generalization of the Chomsky-Schützenberger theorem:

Theorem 30. *Let $L \subseteq \mathbb{T}_{\Sigma}^m$. The following are equivalent:*

- (i) L is simple context-free.
- (ii) There exist finite alphabets $\mathcal{Y}, \mathcal{Y}'$, a projection $\pi: \mathcal{Y}' \rightarrow \mathcal{Y}$, and a local set $R \subseteq \mathbb{T}_{\mathcal{Y}'}^m$, such that $L = \mathbf{del}_{m, \tilde{\mathcal{Y}}-\mathcal{Y}}(\tilde{\pi}(R \cap DT_{\mathcal{Y}'}^m))$.

Proof. (ii) \Rightarrow (i). Suppose $R \subseteq \mathbb{T}_{\mathcal{Y}'}^m$ is a local set. Clearly, $R \subseteq \mathbb{T}_{\mathcal{Y}'_q, p}^m$ for some p, q . So $R \cap DT_{\mathcal{Y}'}^m = R \cap DT_{\mathcal{Y}'}^m \cap \mathbb{T}_{\mathcal{Y}'_q, p}^m$. By Lemma 28, $DT_{\mathcal{Y}'}^m \cap \mathbb{T}_{\mathcal{Y}'_q, p}^m$ is simple context-free. It then follows by Lemma 29 that $L = \mathbf{del}_{m, \tilde{\mathcal{Y}}-\mathcal{Y}}(\tilde{\pi}(R \cap DT_{\mathcal{Y}'}^m)) = \mathbf{del}_{m, \tilde{\mathcal{Y}}-\mathcal{Y}}(\tilde{\pi}(R \cap DT_{\mathcal{Y}'}^m \cap \mathbb{T}_{\mathcal{Y}'_q, p}^m))$ is simple context-free.

(i) \Rightarrow (ii). Let $K \subseteq \mathbb{T}_{\mathcal{Y}, \mathbf{x}}^{m+1}$ be a local set such that $L = \mathbf{y}_m(K)$. By Lemma 27, there exist a projection $\pi: \mathcal{Y}' \rightarrow \mathcal{Y}$, and a local set $R \subseteq \mathbb{T}_{\mathcal{Y}'}^m$ such that $\mathbf{enc}_m(K) = \tilde{\pi}(R \cap DT_{\mathcal{Y}'}^m)$. So

$$\begin{aligned} L &= \mathbf{y}_m(K) \\ &= \mathbf{del}_{m, \tilde{\mathcal{Y}}-\mathcal{Y}}(\mathbf{enc}_m(K)) \\ &= \mathbf{del}_{m, \tilde{\mathcal{Y}}-\mathcal{Y}}(\tilde{\pi}(R \cap DT_{\mathcal{Y}'}^m)). \end{aligned} \quad \square$$

As was the case with the original Chomsky-Schützenberger Theorem, in the proof of Theorem 30, $\mathbf{enc}_m^{-1} \circ \tilde{\pi}$ is a bijection from $R \cap DT_{\mathcal{Y}'}^m$ to K . (See the second statement in Lemma 27.)

9 A Chomsky-Schützenberger-Weir Representation Theorem for Simple Context-Free Tree Grammars

We are now going to use Theorem 30 to obtain a generalization of Weir's representation theorem about the string languages of tree-adjoining grammars to the string languages of simple context-free tree grammars. First, we prove a lemma that generally holds of m -dimensional Dyck tree languages.

The following lemma is straightforward.

Lemma 31. *Let Σ' be a finite alphabet and $\pi: \Sigma' \rightarrow \Sigma$ be a projection. If L is a super-local subset of \mathbb{T}_{Σ}^m , then $\pi^{-1}(L)$ is a super-local subset of $\mathbb{T}_{\Sigma'}^m$.*

Lemma 32. *Let $m \geq 2$. For any local set $L \subseteq \mathbb{T}_{\Sigma}^m$, there exist a finite alphabet Σ' , a degree-bounded, super-local $L' \subseteq \mathbb{T}_{\Sigma'}^m$, and a projection $\pi: \Sigma' \rightarrow \Sigma$ that satisfy the following conditions:*

- (i) $\tilde{\pi}(L') \subseteq L$.

(ii) $L \cap DT_{\Sigma}^m = \tilde{\pi}(L' \cap DT_{\Sigma'}^m)$. Moreover, $\tilde{\pi}$ maps $L' \cap DT_{\Sigma'}^m$ bijectively to $L \cap DT_{\Sigma}^m$.

Proof. By Lemma 10, there are a finite alphabet Σ_1 , a super-local subset L_1 of $\mathbb{T}_{\Sigma_1}^m$, and a projection $\pi_1: \Sigma_1 \rightarrow \tilde{\Sigma}$ such that π_1 maps L_1 bijectively to L . Since $\tilde{\pi}_1^{-1}(L \cap DT_{\Sigma}^m)$ is not a subset of an m -dimensional Dyck tree language, we have to relabel some nodes of $\tilde{\pi}_1^{-1}(\mathbf{T})$ for $\mathbf{T} \in L \cap DT_{\Sigma}^m$ to get a set of the form $L' \cap DT_{\Sigma'}^m$.

For $d \in \Sigma$, P a finite prefix-closed subset of \mathbb{P}_{m-1} , and $i \in [0, |P|]$, let

$$\Delta_{d,P,i} = \{ \delta \in \Sigma_1 \mid \pi_1(\delta) = (d, P, i) \}.$$

Define

$$\begin{aligned} \Sigma_2 &= \Sigma_1 - \bigcup_{d,P,i} \Delta_{d,P,i}, \\ \Delta_{d,P} &= \{ (\delta_0, \delta_1, \dots, \delta_{|P|}) \mid \delta_i \in \Delta_{d,P,i} \text{ for } i = 1, \dots, |P| \}, \\ \Delta &= \bigcup_{d,P} \Delta_{d,P}, \\ \Sigma' &= \Sigma_2 \cup \Delta. \end{aligned}$$

Note that Σ' is a finite alphabet. Define a projection $\pi: \Sigma' \rightarrow \Sigma$ by

$$\begin{aligned} \pi(c) &= \pi_1(c) \quad \text{if } c \in \Sigma_2, \\ \pi(\boldsymbol{\delta}) &= d \quad \text{if } \boldsymbol{\delta} \in \Delta_{d,P}. \end{aligned}$$

Then $\tilde{\pi}$ maps m -dimensional trees over $\tilde{\Sigma}'$ to m -dimensional trees over $\tilde{\Sigma}$. Let

$$\Delta' = \{ (\boldsymbol{\delta}, P, i) \in \tilde{\Sigma}' \mid \boldsymbol{\delta} \in \Delta_{d,P}, 0 \leq i \leq |P| \},$$

and define a projection $\pi_2: \Sigma_2 \cup \Delta' \rightarrow \Sigma_1$ by

$$\begin{aligned} \pi_2(c) &= c \quad \text{if } c \in \Sigma_2, \\ \pi_2(((\delta_0, \delta_1, \dots, \delta_{|P|}), P, i)) &= \delta_i \quad \text{if } (\delta_0, \delta_1, \dots, \delta_k) \in \Delta_{d,P} \text{ and } 0 \leq i \leq |P|. \end{aligned}$$

Then for $\mathbf{T} \in \mathbb{T}_{\Sigma_2 \cup \Delta'}^m$,

$$\tilde{\pi}(\mathbf{T}) = \pi_1(\pi_2(\mathbf{T})).$$

Let

$$L' = \tilde{\pi}^{-1}(L) \cap \mathbb{T}_{\Sigma_2 \cup \Delta'}^m.$$

Then

$$\begin{aligned} L' &= \pi_2^{-1}(\pi_1^{-1}(L)) \\ &= \pi_2^{-1}(L_1). \end{aligned}$$

By Lemma 31, L' is a super-local subset of $\mathbb{T}_{\Sigma_2 \cup \Delta'}^m$, and hence of $\mathbb{T}_{\Sigma'}^m$.

Clearly, $\tilde{\pi}(L') \subseteq L$, so (i) holds. Since $\tilde{\pi}(DT_{\Sigma'}^m) \subseteq DT_{\Sigma}^m$ always holds for any projection $\pi: \Sigma' \rightarrow \Sigma$, we also have $\tilde{\pi}(L' \cap DT_{\Sigma'}^m) \subseteq L \cap DT_{\Sigma}^m$.

It remains to show that for each $\mathbf{T} \in L \cap DT_{\Sigma}^m$, there is a unique $\mathbf{T}' \in L' \cap DT_{\Sigma'}^m$ such that $\tilde{\pi}(\mathbf{T}') = \mathbf{T}$.

Let $\mathbf{T} \in L \cap DT_{\Sigma}^m$. We relabel the nodes of \mathbf{T} to turn it into a $\hat{\mathbf{T}} \in DT_{\Sigma'}^m$. Recall that π_1 maps L_1 bijectively to L , so we have $\mathbf{T}_1 = \pi_1^{-1}(\mathbf{T}) \in L_1$. Let $\mathbf{V} = (V, \ell^{\mathbf{V}}) = \mathbf{enc}_m^{-1}(\mathbf{T})$. Define $\hat{\mathbf{V}} = (V, \ell^{\hat{\mathbf{V}}})$ by

$$\ell^{\hat{\mathbf{V}}}(v) = \begin{cases} \ell^{\mathbf{T}_1}(f_{\mathbf{enc}_m}^{\mathbf{V}}(v)) & \text{if } v \in V - \text{dom}(\prec_{m+1}^V) \text{ and } \ell^{\mathbf{V}}(v) \neq \mathbf{x}, \\ \mathbf{x} & \text{if } \ell^{\mathbf{V}}(v) = \mathbf{x}, \\ (\delta_0, \delta_1, \dots, \delta_k) & \text{if } v \in \text{dom}(\prec_{m+1}^V), |C_m^V(v)| = k, \\ & \delta_0 = \ell^{\mathbf{T}_1}(f_{\mathbf{enc}_m}^{\mathbf{V}}(v)), \text{ and} \\ & v \blacktriangleleft_{m+1,i}^{\mathbf{V}} v_i, \delta_i = \ell^{\mathbf{T}_1}(f_{\mathbf{enc}_m}^{\mathbf{V}}(v_i)) \text{ for } i = 1, \dots, k. \end{cases}$$

Let

$$\hat{\mathbf{T}} = \mathbf{enc}_m(\hat{\mathbf{V}}).$$

If $v \in V - \text{dom}(\prec_{m+1}^V)$ and $\ell^{\hat{\mathbf{V}}}(v) \neq \mathbf{x}$, it is easy to see that $\ell^{\hat{\mathbf{V}}}(v) = \ell^{\hat{\mathbf{T}}}(f_{\mathbf{enc}_m}^{\mathbf{V}}(v)) \in \Sigma_2 \subseteq \Sigma'$. Let $v \in \text{dom}(\prec_{m+1}^V)$, $P = C_m^V(v)$, $k = |P|$, and $v \blacktriangleleft_{m+1,i}^{\mathbf{V}} v_i$ for $i = 1, \dots, k$. Let $(\delta_0, \delta_1, \dots, \delta_k) = \ell^{\hat{\mathbf{V}}}(v)$ and $d = \ell^{\mathbf{V}}(v)$. Then

$$\pi_1(\delta_0) = \ell^{\mathbf{T}}(f_{\mathbf{enc}_m}^{\mathbf{V}}(v)) = (d, P, 0)$$

and for $i = 1, \dots, k$,

$$\pi_1(\delta_i) = \ell^{\mathbf{T}}(f_{\mathbf{enc}_m}^{\mathbf{V}}(v_i)) = (d, P, i).$$

This implies that $(\delta_0, \delta_1, \dots, \delta_k) = \ell^{\hat{\mathbf{V}}}(v) \in \Delta_{d,P} \subseteq \Delta \subseteq \Sigma'$. We have $\ell^{\hat{\mathbf{T}}}(f_{\mathbf{enc}_m}^{\mathbf{V}}(v)) = ((\delta_0, \delta_1, \dots, \delta_k), P, 0) \in \Delta'$ and for $i = 1, \dots, k$, $\ell^{\hat{\mathbf{T}}}(f_{\mathbf{enc}_m}^{\mathbf{V}}(v_i)) = ((\delta_0, \delta_1, \dots, \delta_k), P, i) \in \Delta'$. Therefore, $\hat{\mathbf{V}} \in \mathbb{T}_{\Sigma', \mathbf{x}}^{m+1}$ and $\hat{\mathbf{T}} \in DT_{\Sigma'}^m \cap \mathbb{T}_{\Sigma_2 \cup \Delta'}^m$. It is also easy to see that $\pi_2(\hat{\mathbf{T}}) = \mathbf{T}_1$, so $\hat{\mathbf{T}} \in \pi_2^{-1}(L_1) = L'$. We have shown $\hat{\mathbf{T}} \in L' \cap DT_{\Sigma'}^m$. Since $\pi_2(\hat{\mathbf{T}}) = \mathbf{T}_1$, $\tilde{\pi}(\hat{\mathbf{T}}) = \pi_1(\pi_2(\hat{\mathbf{T}})) = \pi_1(\mathbf{T}_1) = \mathbf{T}$.

Now to show uniqueness, suppose $\mathbf{T}' \in L' \cap DT_{\Sigma'}^m$, and $\mathbf{T} = \tilde{\pi}(\mathbf{T}')$. Let $\mathbf{T}_1 = \pi_1^{-1}(\mathbf{T})$. Then we have $\mathbf{T}_1 = \pi_2(\mathbf{T}')$. We prove $\mathbf{T}' = \hat{\mathbf{T}}$. Let $\mathbf{V}' = (V, \ell^{\mathbf{V}'}) = \mathbf{enc}_m^{-1}(\mathbf{T}')$ and $\mathbf{V} = (V, \ell^{\mathbf{V}}) = \pi(\mathbf{V}')$. Then it is clear that $\mathbf{enc}_m(\mathbf{V}) = \mathbf{T}$. So it suffices to prove $\mathbf{V}' = \hat{\mathbf{V}}$. Let $v \in V$. If $\ell^{\mathbf{V}}(v) = \mathbf{x}$, then clearly, $\ell^{\mathbf{V}'}(v) = \mathbf{x} = \ell^{\hat{\mathbf{V}}}(v)$. If $v \in V - \text{dom}(\prec_{m+1}^V)$ and $\ell^{\mathbf{V}}(v) \neq \mathbf{x}$, then $\ell^{\mathbf{V}'}(v) = \ell^{\mathbf{T}'}(f_{\mathbf{enc}_m}^{\mathbf{V}'}(v)) \in \Sigma_2$, since $\mathbf{V}' \in \mathbb{T}_{\Sigma', \mathbf{x}}^{m+1}$ and $\mathbf{T}' \in L' \subseteq \mathbb{T}_{\Sigma_2 \cup \Delta'}^m$. So $\ell^{\mathbf{V}'}(v) = \ell^{\mathbf{T}'}(f_{\mathbf{enc}_m}^{\mathbf{V}'}(v)) = \pi_2(\ell^{\mathbf{T}'}(f_{\mathbf{enc}_m}^{\mathbf{V}'}(v))) = \ell^{\mathbf{T}_1}(f_{\mathbf{enc}_m}^{\mathbf{V}}(v)) = \ell^{\hat{\mathbf{V}}}(v)$. If $v \in \text{dom}(\prec_{m+1}^V)$ and $C_m^V(v) = P$, then $\ell^{\mathbf{T}'}(f_{\mathbf{enc}_m}^{\mathbf{V}'}(v)) = (\ell^{\mathbf{V}'}(v), P, 0) \in \Delta'$, so $\ell^{\mathbf{V}'}(v) = (\delta_0, \delta_1, \dots, \delta_k)$, where $k = |P|$ and $(\delta_0, \delta_1, \dots, \delta_k) \in \Delta_{d,P}$ for some d . For $i = 1, \dots, k$, let v_i be such that $v \blacktriangleleft_{m+1,i}^{\mathbf{V}'} v_i$, or, equivalently, $v \blacktriangleleft_{m+1,i}^{\mathbf{V}} v_i$. Then for $i = 1, \dots, k$, $\delta_i = \pi_2((\delta_0, \delta_1, \dots, \delta_k), P, i) = \pi_2(\ell^{\mathbf{T}'}(f_{\mathbf{enc}_m}^{\mathbf{V}'}(v_i))) = \ell^{\mathbf{T}_1}(f_{\mathbf{enc}_m}^{\mathbf{V}}(v_i))$. We also have $\delta_0 = \pi_2((\delta_0, \delta_1, \dots, \delta_k), P, 0) = \pi_2(\ell^{\mathbf{T}'}(f_{\mathbf{enc}_m}^{\mathbf{V}'}(v))) = \ell^{\mathbf{T}_1}(f_{\mathbf{enc}_m}^{\mathbf{V}}(v))$. So $\ell^{\hat{\mathbf{V}}}(v) = (\delta_0, \delta_1, \dots, \delta_k) = \ell^{\mathbf{V}'}(v)$. \square

Next we prove a lemma about $DT_{\mathcal{Y}}^2$. Recall that there was an implicit dependence on the dimension m in the definition of $\tilde{\mathcal{Y}}$, which is the alphabet of the language $DT_{\mathcal{Y}}^m$; when a symbol of the form (c, P, i) is in $\tilde{\mathcal{Y}}$, it is assumed that P is a finite, possibly empty, prefix-closed subset of \mathbb{P}_{m-1} . In what follows, we assume that the alphabet $\tilde{\mathcal{Y}}$ is defined from \mathcal{Y} with respect to dimension $m = 2$, so that $(c, P, i) \in \tilde{\mathcal{Y}}$ implies $P = \{\varepsilon, 1, \dots, 1^{k-1}\}$ for some $k \geq 0$. We abbreviate (c, P, i) by (c, k, i) , where $|P| = k$. Under this convention, $\tilde{\mathcal{Y}}_q = \mathcal{Y} \cup \{(c, k, i) \mid c \in \mathcal{Y}, 0 \leq k \leq q, 0 \leq i \leq k\}$.

Recall that for any alphabet Σ , the alphabet Γ_{Σ} consists of symbols of the form \llbracket_c or \rrbracket_c with $c \in \Sigma$.

We let $\mathbb{T}_{\mathcal{Y}, \tilde{\mathcal{Y}}-\mathcal{Y}}$ stand for the set of trees $\mathbf{T} \in \mathbb{T}_{\tilde{\mathcal{Y}}}$ such that for every $v \in T$, $\ell^{\mathbf{T}}(v) \in \tilde{\mathcal{Y}} - \mathcal{Y}$ implies $|C_2^{\mathbf{T}}(v)| = 1$. (In other words, $\tilde{\mathcal{Y}} - \mathcal{Y}$ is regarded as a ranked alphabet all of whose symbols have rank 1.)

Lemma 33. *Let $\eta: \Gamma_{\tilde{\mathcal{Y}}}^* \rightarrow \Gamma_{\mathcal{Y}}^*$ be the alphabetic homomorphism defined as follows:*

$$\begin{aligned} \eta(\llbracket_c) &= \varepsilon, \\ \eta(\rrbracket_c) &= \varepsilon \quad \text{for } c \in \mathcal{Y}, \\ \eta(\llbracket_{(c,k,0)}) &= \llbracket_{(c,k,0)}, \\ \eta(\rrbracket_{(c,k,0)}) &= \rrbracket_{(c,k,k)}, \\ \eta(\llbracket_{(c,k,i)}) &= \rrbracket_{(c,k,i-1)}, \\ \eta(\rrbracket_{(c,k,i)}) &= \llbracket_{(c,k,i)} \quad \text{for } 1 \leq i \leq k. \end{aligned}$$

Then

$$DT_{\mathcal{Y}}^2 = \mathbb{T}_{\mathcal{Y}, \tilde{\mathcal{Y}}-\mathcal{Y}} \cap \mathbf{enc}^{-1}(\eta^{-1}(D_{\tilde{\mathcal{Y}}}).)$$

(Here, \mathbf{enc} is the standard encoding function defined on ordinary 2-dimensional trees.)

Proof. (\subseteq). Suppose $\mathbf{T} \in DT_{\mathcal{Y}}^2$. Clearly, $\mathbf{T} \in \mathbb{T}_{\mathcal{Y}, \tilde{\mathcal{Y}}-\mathcal{Y}}$, so it suffices to prove $\eta(\mathbf{enc}(\mathbf{T})) \in D_{\tilde{\mathcal{Y}}}$. This is proved by induction on the length of the reduction $\mathbf{T} \rightsquigarrow^* \mathbf{T}' \in \mathbb{T}_{\mathcal{Y}}$. If $\mathbf{T} \in \mathbb{T}_{\mathcal{Y}}$, then clearly, $\eta(\mathbf{enc}(\mathbf{T})) = \varepsilon \in D_{\tilde{\mathcal{Y}}}$. Suppose $\mathbf{T} \rightsquigarrow \mathbf{T}'' \rightsquigarrow^* \mathbf{T}' \in \mathbb{T}_{\mathcal{Y}}$. Then $\mathbf{T} = \mathbf{U}[(c, k, 0) \text{ } _{-2} \mathbf{T}_0[(c, k, 1) \text{ } _{-2} \mathbf{T}_1, \dots, (c, k, k) \text{ } _{-2} \mathbf{T}_k]]$ and $\mathbf{T}' = \mathbf{U}[\mathbf{T}_0[\mathbf{T}_1, \dots, \mathbf{T}_k]]$ for some $c \in \mathcal{Y}$, $k \geq 0$, $\mathbf{U} \in \mathbb{T}_{\tilde{\mathcal{Y}}}(1)$, $\mathbf{T}_0 \in \mathbb{T}_{\mathcal{Y}}(k)$, and $\mathbf{T}_i \in \mathbb{T}_{\tilde{\mathcal{Y}}}(i)$ ($i = 1, \dots, k$). Then $\eta(\mathbf{enc}(\mathbf{T})) = z_1 \llbracket_{(c,k,0)} \rrbracket_{(c,k,0)} y_1 \llbracket_{(c,k,1)} \rrbracket_{(c,k,1)} y_2 \dots y_k \llbracket_{(c,k,k)} \rrbracket_{(c,k,k)} z_2$ and $\eta(\mathbf{enc}(\mathbf{T}'')) = z_1 y_1 y_2 \dots y_k z_2$ for some $z_1, z_2, y_1, y_2, \dots, y_k \in (\Gamma_{\tilde{\mathcal{Y}}-\mathcal{Y}})^*$. By induction hypothesis, $\eta(\mathbf{enc}(\mathbf{T}'')) \in D_{\tilde{\mathcal{Y}}}$, and this easily implies $\eta(\mathbf{enc}(\mathbf{T})) \in D_{\tilde{\mathcal{Y}}}$.

(\supseteq). Suppose $\mathbf{T} \in \mathbb{T}_{\mathcal{Y}, \tilde{\mathcal{Y}}-\mathcal{Y}}$ and $\eta(\mathbf{enc}(\mathbf{T})) \in D_{\tilde{\mathcal{Y}}}$. If $\mathbf{T} \in \mathbb{T}_{\mathcal{Y}}$, then $\mathbf{T} \in DT_{\mathcal{Y}}^2$. Suppose that $\mathbf{T} \notin \mathbb{T}_{\mathcal{Y}}$. First, we claim that \mathbf{T} must have a node labeled by $(c, k, 0)$ for some $c \in \mathcal{Y}$ and $k \geq 0$. For, if \mathbf{T} has a node v such that $\ell^{\mathbf{T}}(v) = (c, k, i)$ for some $c \in \mathcal{Y}$, $k \geq 1$, and $i \geq 1$, then $\mathbf{enc}(\mathbf{T}) = x_1 \llbracket_{(c,k,i)} x_2 \rrbracket_{(c,k,i)} x_3$ and $\eta(\mathbf{enc}(\mathbf{T})) = \eta(x_1) \rrbracket_{(c,k,i-1)} \eta(x_2) \llbracket_{(c,k,i)} \eta(x_3)$ for some $x_1, x_2, x_3 \in \Gamma_{\tilde{\mathcal{Y}}}^*$. Since $\eta(\mathbf{enc}(\mathbf{T})) \in D_{\tilde{\mathcal{Y}}}$, $\llbracket_{(c,k,i-1)}$ must occur in $\eta(x_1)$. If $i - 1 = 0$, this implies that $\llbracket_{(c,k,0)}$ occurs in x_1 and it follows that \mathbf{T} has a node labeled by $(c, k, 0)$.

Otherwise, $\mathbb{J}_{(c,k,i-1)}$ occurs in x_1 , and it follows that \mathbf{T} has a node labeled by $(c, k, i - 1)$. Repeating this reasoning, we see that \mathbf{T} must have a node labeled by $(c, k, 0)$.

We show that $\mathbf{T} \in DT_{\mathcal{Y}}^2$ by induction on the number of nodes of \mathbf{T} that are labeled by a symbol of the form $(c, k, 0)$. Let v be a node of \mathbf{T} labeled by $(c, k, 0)$ such that no node v' with $v \prec_2^T v'$ is labeled by a symbol of the form $(d, l, 0)$. Then $\mathbf{enc}(\mathbf{T}) = x_1 \mathbb{J}_{(c,k,0)} y \mathbb{J}_{(c,k,0)} x_2$ for some $x_1, x_2 \in \Gamma_{\mathcal{Y}}^*$ and $y \in \mathbf{enc}(\mathbb{T}_{\mathcal{Y}, \tilde{\mathcal{Y}}-\mathcal{Y}}^2)$. (Note that $|C_2^T(v)| = 1$.)

Case 1. $k = 0$. We show $y \in D'_{\mathcal{Y}}$, i.e., the subtree \mathbf{T}_0 of \mathbf{T} rooted at $v \cdot 2$ is in $\mathbb{T}_{\mathcal{Y}}$. Suppose otherwise, and take the alphabetically first node of \mathbf{T}_0 labeled by some $(d, l, j) \in \tilde{\mathcal{Y}} - \mathcal{Y}$. Then $y = y' \mathbb{J}_{(d,l,j)} y''$ for some $y' \in \Gamma_{\mathcal{Y}}^*$ and $y'' \in \Gamma_{\mathcal{Y}}^*$. By our assumption about v , $j \geq 1$ and $l \geq 1$. We have $\eta(\mathbf{enc}(\mathbf{T})) = \eta(x_1) \mathbb{J}_{(c,0,0)} \eta(y) \mathbb{J}_{(c,0,0)} \eta(x_2) = \eta(x_1) \mathbb{J}_{(c,0,0)} \mathbb{J}_{(d,l,j-1)} \eta(y'') \mathbb{J}_{(c,0,0)} \eta(x_2)$, which contradicts the assumption that $\eta(\mathbf{enc}(\mathbf{T})) \in D_{\tilde{\mathcal{Y}}}$. So \mathbf{T}_0 is in $\mathbb{T}_{\mathcal{Y}}$, and we can write $\mathbf{T} = \mathbf{U}[(c, 0, 0) \cdot_2 \mathbf{T}_0]$. So $\mathbf{T} \rightsquigarrow \mathbf{U}[\mathbf{T}_0] \in \mathbb{T}_{\mathcal{Y}, \tilde{\mathcal{Y}}-\mathcal{Y}}$. We have $\eta(\mathbf{enc}(\mathbf{T})) = \eta(x_1) \mathbb{J}_{(c,0,0)} \mathbb{J}_{(c,0,0)} \eta(x_2)$ and $\eta(\mathbf{enc}(\mathbf{U}[\mathbf{T}_0])) = \eta(x_1)\eta(x_2)$. Since $\eta(\mathbf{enc}(\mathbf{T})) \in D_{\tilde{\mathcal{Y}}}$, $\eta(x_1)\eta(x_2) \in D_{\tilde{\mathcal{Y}}}$ as well, and $\mathbf{U}[\mathbf{T}_0] \in DT_{\mathcal{Y}}^2$ by induction hypothesis. Since $\mathbf{T} \rightsquigarrow \mathbf{U}[\mathbf{T}_0]$, we conclude $\mathbf{T} \in DT_{\mathcal{Y}}^2$.

Case 2. $k \geq 1$. We show

$$y = z_0 \mathbb{J}_{(c,k,1)} y_1 \mathbb{J}_{(c,k,1)} z_1 \mathbb{J}_{(c,k,2)} y_2 \mathbb{J}_{(c,k,2)} \cdots z_{k-1} \mathbb{J}_{(c,k,k)} y_k \mathbb{J}_{(c,k,k)} z_k$$

for some $z_0, z_1, \dots, z_k \in \Gamma_{\mathcal{Y}}^*$ and $y_1, y_2, \dots, y_k \in \mathbf{enc}(\mathbb{T}_{\mathcal{Y}, \tilde{\mathcal{Y}}-\mathcal{Y}})$. First, we show by induction that the following condition holds for $i = 0, \dots, k$:

$$y \text{ has a prefix of the form } z_0 \mathbb{J}_{(c,k,1)} y_1 \mathbb{J}_{(c,k,1)} \cdots z_{i-1} \mathbb{J}_{(c,k,i)} y_i \mathbb{J}_{(c,k,i)}. \quad (12)$$

The case of $i = 0$ is trivial. Suppose we have shown (12) for $i < k$, i.e., $y = z_0 \mathbb{J}_{(c,k,1)} y_1 \mathbb{J}_{(c,k,1)} \cdots z_{i-1} \mathbb{J}_{(c,k,i)} y_i \mathbb{J}_{(c,k,i)} y'$ with $z_0, \dots, z_{i-1} \in \Gamma_{\mathcal{Y}}^*$ and $y_1, \dots, y_i \in \mathbf{enc}(\mathbb{T}_{\mathcal{Y}, \tilde{\mathcal{Y}}-\mathcal{Y}})$. Since

$$\begin{aligned} \eta(\mathbf{enc}(\mathbf{T})) &= \eta(x_1) \mathbb{J}_{(c,k,0)} \eta(y) \mathbb{J}_{(c,k,k)} \eta(x_2) \\ &= \eta(x_1) \mathbb{J}_{(c,k,0)} \mathbb{J}_{(c,k,0)} \eta(y_1) \mathbb{J}_{(c,k,1)} \cdots \mathbb{J}_{(c,k,i-1)} \eta(y_i) \mathbb{J}_{(c,k,i)} \eta(y') \mathbb{J}_{(c,k,k)}, \end{aligned}$$

we must have $\eta(y') \neq \varepsilon$. Let z_i be the longest prefix of y' in $\Gamma_{\mathcal{Y}}^*$. Since $z_0, \dots, z_{i-1} \in \Gamma_{\mathcal{Y}}^*$ and y_1, \dots, y_i and y are all in $\mathbf{enc}(\mathbb{T}_{\mathcal{Y}, \tilde{\mathcal{Y}}-\mathcal{Y}})$, it is easy to see that $y' = z_i \mathbb{J}_{(d,l,j)} y_{i+1} \mathbb{J}_{(d,l,j)} y''$ for some $y_{i+1} \in \mathbf{enc}(\mathbb{T}_{\mathcal{Y}, \tilde{\mathcal{Y}}-\mathcal{Y}})$ and $l, j \geq 1$. We have $\eta(y') = \mathbb{J}_{(d,l,j-1)} \eta(y_{i+1}) \mathbb{J}_{(d,l,j)} \eta(y'')$, and so

$$\begin{aligned} \eta(\mathbf{enc}(\mathbf{T})) &= \\ \eta(x_1) \mathbb{J}_{(c,k,0)} \mathbb{J}_{(c,k,0)} \eta(y_1) \mathbb{J}_{(c,k,1)} \cdots \mathbb{J}_{(c,k,i-1)} \eta(y_i) \mathbb{J}_{(c,k,i)} \mathbb{J}_{(d,l,j-1)} \eta(y_{i+1}) \mathbb{J}_{(d,l,j)} \eta(y'') \mathbb{J}_{(c,k,k)}, \end{aligned}$$

which implies $d = c$, $l = k$, and $j = i + 1$. This shows that (12) holds with $i + 1$ in place of i . By induction, (12) holds with $i = k$.

We have

$$y = z_0 \llbracket_{(c,k,1)} y_1 \rrbracket_{(c,k,1)} \cdots z_{k-1} \llbracket_{(c,k,k)} y_k \rrbracket_{(c,k,k)} y'$$

with $z_0, \dots, z_{k-1} \in \Gamma_{\mathcal{Y}}^*$ and $y_1, \dots, y_k \in \mathbf{enc}(\mathbb{T}_{\mathcal{Y}, \tilde{\mathcal{Y}}-\mathcal{Y}})$. We show that $y' \in \Gamma_{\mathcal{Y}}^*$. Suppose otherwise. Then we must have $y' = z_k \llbracket_{(d,l,j)} y_{k+1} \rrbracket_{(d,l,j)} y''$ for some $d \in \mathcal{Y}$, $j, l \geq 1$, $z_k \in \Gamma_{\mathcal{Y}}^*$, and $y_{k+1} \in \mathbf{enc}(\mathbb{T}_{\mathcal{Y}, \tilde{\mathcal{Y}}-\mathcal{Y}})$. Then $\eta(\mathbf{enc}(\mathbf{T}))$ contains as a substring $\llbracket_{(c,k,k)} \rrbracket_{(d,l,j-1)}$, which is a contradiction since $\eta(\mathbf{enc}(\mathbf{T})) \in D_{\tilde{\mathcal{Y}}}$ and $j-1 < l$. Therefore,

$$y = z_0 \llbracket_{(c,k,1)} y_1 \rrbracket_{(c,k,1)} \cdots z_{k-1} \llbracket_{(c,k,k)} y_k \rrbracket_{(c,k,k)} z_k$$

with $z_0, \dots, z_k \in \Gamma_{\mathcal{Y}}^*$ and $y_1, \dots, y_k \in \mathbf{enc}(\mathbb{T}_{\mathcal{Y}, \tilde{\mathcal{Y}}-\mathcal{Y}})$. This means that \mathbf{T} is of the form

$$\mathbf{T} = \mathbf{U}[(c, k, 0) \text{--}_2 \mathbf{T}_0[(c, k, 1) \text{--}_2 \mathbf{T}_1, \dots, (c, k, k) \text{--}_2 \mathbf{T}_k]],$$

where $\mathbf{T}_0 \in \mathbb{T}_{\mathcal{Y}}(k)$. So

$$\mathbf{T} \rightsquigarrow \mathbf{T}' = \mathbf{U}[\mathbf{T}_0[\mathbf{T}_1, \dots, \mathbf{T}_k]].$$

Clearly, $\mathbf{T}' \in \mathbb{T}_{\mathcal{Y}, \tilde{\mathcal{Y}}-\mathcal{Y}}$, and

$$\eta(\mathbf{enc}(\mathbf{T}')) = \eta(x_1)\eta(y_1) \cdots \eta(y_k)\eta(x_2).$$

Since

$$\eta(\mathbf{enc}(\mathbf{T})) = \eta(x_1) \llbracket_{(c,k,0)} \rrbracket_{(c,k,0)} \eta(y_1) \llbracket_{(c,k,1)} \rrbracket_{(c,k,1)} \eta(y_2) \cdots \llbracket_{(c,k,k-1)} \rrbracket_{(c,k,k-1)} \eta(y_k) \llbracket_{(c,k,k)} \rrbracket_{(c,k,k)} \eta(x_2)$$

is in $D_{\tilde{\mathcal{Y}}}$, it follows that $\eta(\mathbf{enc}(\mathbf{T}'))$ is in $D_{\tilde{\mathcal{Y}}}$ as well, and the induction hypothesis gives $\mathbf{T}' \in DT_{\mathcal{Y}}^2$. Since $\mathbf{T} \rightsquigarrow \mathbf{T}'$, we conclude $\mathbf{T} \in DT_{\mathcal{Y}}^2$. \square

Lemma 34. *If $L \subseteq \mathbb{T}_{\Sigma, \mathbf{x}}^3$ is a local set, then there exist a finite alphabet \mathcal{Y} , a projection $\pi: \mathcal{Y} \rightarrow \Sigma$, and a local set $R \subseteq \Gamma_{\mathcal{Y}}^+$ such that*

$$\mathbf{enc}(\mathbf{enc}_2(L)) = \widehat{\pi}(R \cap D_{\tilde{\mathcal{Y}}} \cap \eta^{-1}(D_{\tilde{\mathcal{Y}}}),$$

where η is the alphabetic homomorphism defined in Lemma 33. Moreover, $\mathbf{enc}_2^{-1} \circ \mathbf{enc}^{-1} \circ \widehat{\pi}$ maps $R \cap D_{\tilde{\mathcal{Y}}} \cap \eta^{-1}(D_{\tilde{\mathcal{Y}}})$ bijectively to L .

Proof. By Lemma 27, there exist a finite alphabet \mathcal{Y}_1 , a projection $\pi_1: \mathcal{Y}_1 \rightarrow \Sigma$, and a local set $L_1 \subseteq \mathbb{T}_{\mathcal{Y}_1}^2$ such that $\mathbf{enc}_2(L) = \widehat{\pi}_1(L_1 \cap DT_{\mathcal{Y}_1}^2)$ and $\widehat{\pi}_1$ is a bijection from $L_1 \cap DT_{\mathcal{Y}_1}^2$ to $\mathbf{enc}_2(L)$. We may assume $L_1 \subseteq \mathbb{T}_{\mathcal{Y}_1, \tilde{\mathcal{Y}}_1-\mathcal{Y}_1}$. By Lemma 32, there exist a finite alphabet \mathcal{Y}_2 , a projection $\pi_2: \mathcal{Y}_2 \rightarrow \mathcal{Y}_1$, and a super-local set $L_2 \subseteq \mathbb{T}_{\mathcal{Y}_2}^2$ such that $\widehat{\pi}_2(L_2) \subseteq L_1$ and $\widehat{\pi}_2(L_2 \cap DT_{\mathcal{Y}_2}^2) = L_1 \cap DT_{\mathcal{Y}_1}^2$. Since $L_1 \subseteq \mathbb{T}_{\mathcal{Y}_1, \tilde{\mathcal{Y}}_1-\mathcal{Y}_1}$, it follows that $L_2 \subseteq \mathbb{T}_{\mathcal{Y}_2, \tilde{\mathcal{Y}}_2-\mathcal{Y}_2}$. We have

$$\mathbf{enc}(\mathbf{enc}_2(L)) = \mathbf{enc}(\widehat{\pi}_1(L_1 \cap DT_{\mathcal{Y}_1}^2))$$

$$\begin{aligned}
&= \mathbf{enc}(\widehat{\pi}_1(\widehat{\pi}_2(L_2 \cap DT_{\mathcal{Y}_2}^2))) \\
&= \widehat{\pi}_1(\widehat{\pi}_2(\mathbf{enc}(L_2 \cap DT_{\mathcal{Y}_2}^2))) \\
&= \widehat{\pi}_1(\widehat{\pi}_2(\mathbf{enc}(L_2) \cap \mathbf{enc}(DT_{\mathcal{Y}_2}^2))), \tag{13}
\end{aligned}$$

since \mathbf{enc} is injective. By Lemma 33,

$$\mathbf{enc}(DT_{\mathcal{Y}_2}^2) = \mathbf{enc}(\mathbb{T}_{\mathcal{Y}_2, \widetilde{\mathcal{Y}}_2 - \mathcal{Y}_2}) \cap \eta_2^{-1}(D_{\widetilde{\mathcal{Y}}_2}),$$

where $\eta_2: \Gamma_{\widetilde{\mathcal{Y}}_2}^* \rightarrow \Gamma_{\mathcal{Y}_2}^*$ is an alphabetic homomorphism defined like η . Since $L_2 \subseteq \mathbb{T}_{\mathcal{Y}_2, \widetilde{\mathcal{Y}}_2 - \mathcal{Y}_2}$,

$$\mathbf{enc}(L_2) \cap \mathbf{enc}(DT_{\mathcal{Y}_2}^2) = \mathbf{enc}(L_2) \cap \eta_2^{-1}(D_{\widetilde{\mathcal{Y}}_2}).$$

By Lemma 2, $\mathbf{enc}(L_2) = R_2 \cap D'_{\widetilde{\mathcal{Y}}_2}$ for some local set $R_2 \subseteq \Gamma_{\widetilde{\mathcal{Y}}_2}^+$. So we have

$$\mathbf{enc}(L_2) \cap \mathbf{enc}(DT_{\mathcal{Y}_2}^2) = R_2 \cap D'_{\widetilde{\mathcal{Y}}_2} \cap \eta_2^{-1}(D_{\widetilde{\mathcal{Y}}_2}). \tag{14}$$

Given (13) and (14), all we need is to turn $R_2 \cap D'_{\widetilde{\mathcal{Y}}_2} \cap \eta_2^{-1}(D_{\widetilde{\mathcal{Y}}_2})$ into the form $\widehat{\pi}_3(R \cap D_{\widetilde{\mathcal{Y}}} \cap \eta^{-1}(D_{\widetilde{\mathcal{Y}}}))$. For this, we can use a method similar to the one we used in the proof of Lemma 6. Let

$$\mathcal{Y} = \mathcal{Y}_2 \cup \{ \bar{c} \mid c \in \mathcal{Y}_2 \}$$

and define $\pi_3: \mathcal{Y} \rightarrow \mathcal{Y}_2$ by

$$\pi_3(c) = c, \quad \pi_3(\bar{c}) = c,$$

for each $c \in \mathcal{Y}_2$. Let

$$\begin{aligned}
\Delta_1 &= \{ \bar{c} \mid c \in \mathcal{Y}_2 \} \cup \{ (\bar{c}, P, 0) \mid (c, P, 0) \in \widetilde{\mathcal{Y}}_2 \}, \\
\Delta_2 &= \widetilde{\mathcal{Y}}_2 \cup \{ (\bar{c}, P, i) \mid (c, P, i) \in \widetilde{\mathcal{Y}}_2, i \geq 1 \}.
\end{aligned}$$

Then Δ_1, Δ_2 is a partition of $\widetilde{\mathcal{Y}}$. Let

$$R = (\{ \llbracket d \mid d \in \Delta_1 \rrbracket \Gamma_{\Delta_2}^* \{ \rrbracket d \mid d \in \Delta_1 \rrbracket \} \cap \widehat{\pi}_3^{-1}(R_2)).$$

Then R is a local subset of $\Gamma_{\widetilde{\mathcal{Y}}}^+$,²⁶ and it is easy to see

$$\widehat{\pi}_3(R \cap D_{\widetilde{\mathcal{Y}}} \cap \eta^{-1}(D_{\widetilde{\mathcal{Y}}})) = R_2 \cap D'_{\widetilde{\mathcal{Y}}_2} \cap \eta_2^{-1}(D_{\widetilde{\mathcal{Y}}_2}), \tag{15}$$

where η^{-1} is as defined in Lemma 33. It is also easy to see that $\widehat{\pi}_3$ maps $R \cap D_{\widetilde{\mathcal{Y}}} \cap \eta^{-1}(D_{\widetilde{\mathcal{Y}}})$ bijectively to $R_2 \cap D'_{\widetilde{\mathcal{Y}}_2} \cap \eta_2^{-1}(D_{\widetilde{\mathcal{Y}}_2})$.

We obtain the statement of the lemma from (13), (14), and (15) by letting $\pi = \pi_3 \circ \pi_2 \circ \pi_1$. \square

²⁶ Although $\Gamma_{\widetilde{\mathcal{Y}}}^+$ is an infinite alphabet, only finitely many symbols in it appear in $\widehat{\pi}_3^{-1}(R_2)$ since R_2 is local.

Recall that when Σ_0 and Σ_1 are disjoint alphabets, $\mathbb{T}_{\Sigma_0}^{\Sigma_1}$ consists of all (ordinary 2-dimensional) trees in $\mathbb{T}_{\Sigma_0 \cup \Sigma_1}$ that are disjointly labeled with Σ_0, Σ_1 .

Lemma 35. *If $L \subseteq \mathbb{T}_{\Sigma, \mathbf{x}}^3$ is a local set, then there exist an alphabet Σ' disjoint from Σ , a projection $\pi: \Sigma \cup \Sigma' \rightarrow \Sigma$, and a local set $L' \subseteq \mathbb{T}_{\Sigma \cup \Sigma', \mathbf{x}}^3$ that satisfy the following conditions:*

- (i) $\pi(c) = c$ for all $c \in \Sigma$.
- (ii) $L = \pi(L')$. Moreover, π maps L' bijectively to L .
- (iii) $\mathbf{y}_2(L') \subseteq \mathbb{T}_{\Sigma'}^{\Sigma}$.
- (iv) $\mathbf{y}_2(L) = \pi(\mathbf{y}_2(L'))$. Moreover, π maps $\mathbf{y}_2(L')$ bijectively to $\mathbf{y}_2(L)$.
- (v) $\mathbf{y}(\mathbf{y}_2(L)) = \mathbf{y}(\mathbf{y}_2(L'))$.

Proof. Let $L = \text{Loc}^3(A, Z, I)$ be a local subset of $\mathbb{T}_{\Sigma, \mathbf{x}}^3$. Let $\Sigma' = \{\bar{c} \mid c \in \Sigma\}$. Define a projection $\pi: \Sigma \cup \Sigma' \rightarrow \Sigma$ by

$$\pi(c) = c, \quad \pi(\bar{c}) = c,$$

for each $c \in \Sigma$. Let

$$Z' = Z \cup \{\bar{c} \mid c \in Z\}$$

$$I' = \{(c, \mathbf{T}') \mid (c, \mathbf{T}) \in I, \mathbf{T} \in \mathbb{T}_{\Sigma}^2(0)\} \cup \{(\bar{c}, \mathbf{T}') \mid (c, \mathbf{T}) \in I, \mathbf{T} \in \mathbb{T}_{\Sigma}^2(n), n \geq 1\},$$

where for $\mathbf{T} = (T, \ell^{\mathbf{T}}) \in \mathbb{T}_{\Sigma \cup \{\mathbf{x}\}}^2$, $\mathbf{T}' = (T, \ell^{\mathbf{T}'})$ is defined by

$$\ell^{\mathbf{T}'}(v) = \begin{cases} \bar{c} & \text{if } \ell^{\mathbf{T}}(v) = c \in \Sigma \text{ and } v \in \text{dom}(\prec_2^T), \\ \ell^{\mathbf{T}}(v) & \text{otherwise.} \end{cases}$$

Note that $\mathbf{T}' \in \mathbb{T}_{\Sigma \cup \Sigma' \cup \{\mathbf{x}\}}^2$. Define a local subset L' of $\mathbb{T}_{\Sigma \cup \Sigma' \cup \{\mathbf{x}\}}^3$ by $L' = \text{Loc}^3(A, Z', I')$. Then it is easy to see that π and L' satisfy the required properties. \square

Lemma 36. *If $L \subseteq \mathbb{T}_{\Sigma, \mathbf{x}}^3$ is a local set, then there exist a finite alphabet Υ , a local set $R \subseteq \Gamma_{\Upsilon}^+$, and an alphabetic homomorphism $h: \Gamma_{\Upsilon}^* \rightarrow \Sigma^*$ such that*

$$\mathbf{y}(\mathbf{y}_2(L)) = h(R \cap D_{\tilde{\Upsilon}} \cap \eta^{-1}(D_{\tilde{\Upsilon}})),$$

where $\tilde{\Upsilon}$ is defined with respect to dimension 2 and η is the alphabetic homomorphism defined in Lemma 33.

Proof. Let $L \subseteq \mathbb{T}_{\Sigma, \mathbf{x}}^3$ be a local set. Applying Lemma 35 to L , we obtain a projection $\pi': \Sigma \cup \Sigma' \rightarrow \Sigma$ and a local set $L' \subseteq \mathbb{T}_{\Sigma \cup \Sigma', \mathbf{x}}^3$ such that π' maps L' bijectively to L , $\mathbf{y}_2(L') \subseteq \mathbb{T}_{\Sigma'}^{\Sigma}$, and $\mathbf{y}(\mathbf{y}_2(L)) = \mathbf{y}(\mathbf{y}_2(L'))$. By Lemma 34,

$$\mathbf{enc}(\mathbf{enc}_2(L')) = \hat{\pi}(R \cap D_{\tilde{\Upsilon}} \cap \eta^{-1}(D_{\tilde{\Upsilon}}))$$

for some finite set \mathcal{T} , projection $\pi: \mathcal{Y} \rightarrow \Sigma \cup \Sigma'$, and local set $R \subseteq \Gamma_{\tilde{\mathcal{Y}}}^+$. We have

$$\begin{aligned} \mathbf{y}(\mathbf{y}_2(L)) &= \mathbf{y}(\mathbf{y}_2(L')) \\ &= h_{\Sigma_0, \Sigma_1}(\mathbf{enc}(\mathbf{del}_{2, \tilde{\mathcal{Y}}-\mathcal{T}}(\mathbf{enc}_2(L')))) \\ &= h_{\Sigma_0, \Sigma_1}(h_{\Gamma_{\mathcal{T}}, \Gamma_{\tilde{\mathcal{Y}}-\mathcal{T}}}(\mathbf{enc}(\mathbf{enc}_2(L')))) \\ &= h_{\Sigma_0, \Sigma_1}(h_{\Gamma_{\mathcal{T}}, \Gamma_{\tilde{\mathcal{Y}}-\mathcal{T}}}(\widehat{\pi}(R \cap D_{\tilde{\mathcal{Y}}} \cap \eta^{-1}(D_{\tilde{\mathcal{Y}}})))), \end{aligned}$$

so the statement of the lemma holds with $h = h_{\Sigma_0, \Sigma_1} \circ h_{\Gamma_{\mathcal{T}}, \Gamma_{\tilde{\mathcal{Y}}-\mathcal{T}}} \circ \widehat{\pi}$. \square

Note that in the above proof, the set $R \cap D_{\tilde{\mathcal{Y}}} \cap \eta^{-1}(D_{\tilde{\mathcal{Y}}})$ is mapped bijectively to L by $\pi' \circ \mathbf{enc}_2^{-1} \circ \mathbf{enc}^{-1} \circ \widehat{\pi}$.

The following lemma is analogous to the corresponding characterization of context-free (string) languages.

Lemma 37. *Let $L \subseteq \mathbb{T}_{\Sigma}$. Then $L \in \text{CFT}_{\text{sp}}(r)$ if and only if there exist a finite alphabet \mathcal{T} and a local set $K \subseteq \mathbb{T}_{\mathcal{T}, \mathbf{x}}^3$ such that*

- (i) $L = \mathbf{y}_2(K)$, and
- (ii) for all $\mathbf{T} \in K$ and all $v \in T$, if $v \in \text{dom}(\prec_3^T)$, then $|C_2^T(v)| \leq r$.

Lemma 38. *For every $L \in \text{CFT}_{\text{sp}}(r)$, there is an $L' \in \text{CFT}_{\text{sp}}(r)$ such that $\mathbf{enc}(L) = \mathbf{y}(L')$.*

Proof. Let $K \subseteq \mathbb{T}_{\mathcal{T}, \mathbf{x}}^3$ be a local set satisfying condition (ii) of Lemma 37. By Lemma 35, we may assume $K = \text{Loc}^3(A, Z, I)$ with $Z \cap \{c \mid (c, \mathbf{T}) \in I\} = \emptyset$. Let

$$\begin{aligned} A' &= A \cup \{\bar{c} \mid c \in A \cap Z\}, \\ Z' &= Z \cup \bigcup \{ \{ \llbracket c, \rrbracket c \} \mid c \in Z \}, \\ I' &= \{ (c, \varphi(\mathbf{T})) \mid (c, \mathbf{T}) \in I \} \cup \{ (\bar{c}, c(\llbracket c, \rrbracket c)) \mid c \in A \cap Z \}. \end{aligned}$$

where for each $c \in A \cap Z$, \bar{c} is a new symbol, and

$$\begin{aligned} \varphi(\mathbf{x}) &= \mathbf{x}, \\ \varphi(c) &= \begin{cases} c(\llbracket c, \rrbracket c) & \text{if } c \in Z, \\ c & \text{otherwise,} \end{cases} \\ \varphi(c(\mathbf{T}_1 \dots \mathbf{T}_n)) &= \begin{cases} c(\llbracket c, \varphi(\mathbf{T}_1) \dots \varphi(\mathbf{T}_n) \rrbracket c) & \text{if } c \in Z, \\ c(\varphi(\mathbf{T}_1) \dots \varphi(\mathbf{T}_n)) & \text{otherwise.} \end{cases} \end{aligned}$$

Then it is easy to see $K' = \text{Loc}^3(A, Z', I')$ also satisfies condition (ii) of Lemma 37, and we have $\mathbf{enc}(\mathbf{y}_2(K)) = \mathbf{y}(\mathbf{y}_2(K'))$. We omit the details. \square

Recall that $\Gamma_n = \{\llbracket_1, \rrbracket_1, \dots, \llbracket_n, \rrbracket_n\}$ and D_n is the Dyck language over Γ_n , where \llbracket_i and \rrbracket_i form a matching pair of brackets for $i = 1, \dots, n$. Define a bijection $g: \Gamma_{qn} \rightarrow \Gamma_{qn}$ by

$$\begin{aligned} g(\llbracket_{qi+1}) &= \llbracket_{qi+1}, & g(\rrbracket_{qi+1}) &= \rrbracket_{qi+q}, \\ g(\llbracket_{qi+j}) &= \rrbracket_{qi+j-1}, & g(\rrbracket_{qi+j}) &= \llbracket_{qi+j}, \end{aligned} \quad (16)$$

for $i = 0, \dots, n-1$ and $j = 2, \dots, q$.

Lemma 39. For $q, n \geq 1$, $D_{qn} \cap g^{-1}(D_{qn}) \in \text{yCFT}_{\text{sp}}(q-1)$.

Proof. By Lemma 38 and the fact that $\text{yCFT}_{\text{sp}}(q-1)$ is a substitution-closed full abstract family of languages [37], it suffices to show that there are some $L \in \text{CFT}_{\text{sp}}(q-1)$ and homomorphism h such that $h(\mathbf{enc}(L)) = D_{qn} \cap g^{-1}(D_{qn})$. Let $m = 2$, $\Sigma = \{c_1, \dots, c_n\}$, $r = q-1$, $p = 2$, and $\mathcal{Y} = \tilde{\Sigma}_{q-1} \cup \{X_0, \dots, X_{q-1}\}$. Lemma 28 gives finite sets $A \subseteq \mathcal{Y}$, $Z = \tilde{\Sigma}_{q-1}$, $I \subseteq \{X_0, \dots, X_{q-1}\} \times \mathbb{T}_{\mathcal{Y}}^2$ such that

$$DT_{\Sigma}^2 \cap \mathbb{T}_{\tilde{\Sigma}_{q-1}, 2}^2 = \mathbf{y}_2(\text{Loc}^3(A, Z, I)).$$

Let $L = DT_{\Sigma}^2 \cap \mathbb{T}_{\tilde{\Sigma}_{q-1}, 2}^2$. Inspection of the proof of Lemma 28 also shows that for all $\mathbf{T} \in \text{Loc}^3(A, Z, I)$ and all $v \in T$, $v \in \text{dom}(\prec_3^T)$ implies $|C_2^T(v)| \leq q-1$. So $L \in \text{CFT}_{\text{sp}}(q-1)$. Let $\Phi = \{(c_i, q-1, j) \mid 1 \leq i \leq n, 0 \leq j \leq q-1\}$. We identify $\Gamma_{\Phi} = \bigcup \{\{\llbracket_d, \rrbracket_d\} \mid d \in \Phi\}$ with Γ_{qn} . Let $h: (\Gamma_{\tilde{\Sigma}_{q-1}}^*)^* \rightarrow (\Gamma_{\tilde{\Sigma}_{q-1}}^*)^*$ be the homomorphism that erases all symbols that are not in Γ_{qn} . Our goal is to show

$$h(\mathbf{enc}(L)) = D_{qn} \cap g^{-1}(D_{qn}).$$

To show $h(\mathbf{enc}(L)) \subseteq D_{qn} \cap g^{-1}(D_{qn})$, suppose $\mathbf{T} \in L$. Since $L \subseteq \mathbb{T}_{\tilde{\Sigma}_{q-1}, 2}^2$, it is clear that

$$h(\mathbf{enc}(\mathbf{T})) \in D_{qn}. \quad (17)$$

Since $L \subseteq DT_{\Sigma}^2$, by Lemma 33,

$$\eta(\mathbf{enc}(\mathbf{T})) \in D_{\tilde{\Sigma}_{q-1}},$$

where $\eta: \Gamma_{\tilde{\Sigma}}^* \rightarrow \Gamma_{\tilde{\Sigma}}^*$ is as defined in Lemma 33, with Σ in place of \mathcal{Y} . So

$$h(\eta(\mathbf{enc}(\mathbf{T}))) \in h(D_{\tilde{\Sigma}_{q-1}}) = D_{qn}.$$

Note that η restricted to Γ_{qn} coincides with g . So we have

$$\begin{aligned} h(\eta(\mathbf{enc}(\mathbf{T}))) &= \eta(h(\mathbf{enc}(\mathbf{T}))) \\ &= g(h(\mathbf{enc}(\mathbf{T}))). \end{aligned}$$

This shows that

$$h(\mathbf{enc}(\mathbf{T})) \in g^{-1}(D_{qn}). \quad (18)$$

By (17) and (18), $h(\mathbf{enc}(L)) \subseteq D_{qn} \cap g^{-1}(D_{qn})$.

Now we show the converse inclusion. Let $s \in D_{qn} \cap g^{-1}(D_{qn})$. Then there is a hedge $\mathbf{T} \in \mathbb{H}_{\Phi}^2$ such that $s = \mathbf{enc}(\mathbf{T})$. We turn \mathbf{T} into a tree $\mathbf{T}' = \varphi(\mathbf{T}) \in \mathbb{T}_{\{c_1\}, \Phi} \cap \mathbb{T}_{\Sigma_{q-1}, 2}^2$, where symbols in Φ are assumed to have rank 1:

$$\begin{aligned}\varphi((c_i, q-1, j)) &= (c_i, q-1, j)(c_1), \\ \varphi((c_i, q-1, j)(\mathbf{T}_1 \dots \mathbf{T}_n)) &= (c_i, q-1, j)(\varphi(\mathbf{T}_1 \dots \mathbf{T}_n)), \\ \varphi(\mathbf{T}_1 \dots \mathbf{T}_n) &= c_1(\varphi(\mathbf{T}_1) \varphi(\mathbf{T}_2 \dots \mathbf{T}_n)) \quad \text{where } n \geq 2.\end{aligned}$$

Then $s = h(\mathbf{enc}(\mathbf{T}'))$. We have

$$\eta(\mathbf{enc}(\mathbf{T}')) = g(s) \in D_{qn} \subseteq D_{\tilde{\Sigma}}.$$

Since $\mathbf{T}' \in \mathbb{T}_{\{c_1\}, \Phi} \subseteq \mathbb{T}_{\Sigma, \tilde{\Sigma}-\Sigma}$, Lemma 33 implies that $\mathbf{T}' \in DT_{\Sigma}^2$. So $\mathbf{T}' \in L$ and $s = h(\mathbf{enc}(\mathbf{T}')) \in h(\mathbf{enc}(L))$. We conclude $D_{qn} \cap g^{-1}(D_{qn}) \subseteq h(\mathbf{enc}(L))$.

We have shown $h(\mathbf{enc}(L)) = D_{qn} \cap g^{-1}(D_{qn})$. \square

Theorem 40. *Let $q \geq 1$ and $M \subseteq \Sigma^*$. The following are equivalent:*

- (i) $M \in \text{yCFT}_{\text{sp}}(q-1)$.
- (ii) *There exist a positive integer n , a local set $R \subseteq \Gamma_{qn}^*$, and an alphabetic homomorphism $h: \Gamma_{qn}^* \rightarrow \Sigma^*$ such that $M = h(R \cap D_{qn} \cap g^{-1}(D_{qn}))$.*

Proof. (ii) \Rightarrow (i). This follows from Lemma 39 and the fact that $\text{yCFT}_{\text{sp}}(q-1)$ is an abstract family of languages.

(i) \Rightarrow (ii). Let $L \in \text{CFT}_{\text{sp}}(q-1)$ be such that $M = \mathbf{y}(L)$. By Lemma 37, $L = \mathbf{y}_2(K)$ and $\mathbf{enc}_2(K) \subseteq \mathbb{T}_{\Psi_{q-1}}^2$ for some finite set Ψ and some local set $K \subseteq \mathbb{T}_{\Psi, \mathbf{x}}^3$. By Lemma 34, $\mathbf{enc}(\mathbf{enc}_2(K)) = \widehat{\pi}(R' \cap D_{\tilde{\Upsilon}} \cap \eta^{-1}(D_{\tilde{\Upsilon}}))$ for some local set $R' \subseteq \Gamma_{\tilde{\Upsilon}}$ and projection $\pi: \Upsilon \rightarrow \Psi$. Since $\mathbf{enc}_2(K) \subseteq \mathbb{T}_{\Psi_{q-1}}^2$, it easily follows that $\mathbf{enc}(\mathbf{enc}_2(K)) = \widehat{\pi}(R'' \cap D_{\tilde{\Upsilon}_q} \cap \eta^{-1}(D_{\tilde{\Upsilon}_q}))$, where $R'' = R' \cap (\Gamma_{\tilde{\Upsilon}_{q-1}})^+$, which is a local subset of $(\Gamma_{\tilde{\Upsilon}_{q-1}})^+$. Using this in the proof of Lemma 36, we easily obtain

$$M = h'(R'' \cap D_{\tilde{\Upsilon}_{q-1}} \cap \eta^{-1}(D_{\tilde{\Upsilon}_{q-1}})), \quad (19)$$

where $h': (\Gamma_{\tilde{\Upsilon}_{q-1}})^* \rightarrow \Sigma^*$ is an alphabetic homomorphism. In order to obtain the statement of the lemma, there are three things we need to fix:

- for $c \in \Upsilon$, η erases $[c$ and $]$,
- when $q \geq 2$, the number of pairs of brackets in the group $[c,]_c$ is $1 < q$, and
- when $k < q-1$, the number of pairs of brackets in the group $[_{(c,k,0)},]_{(c,k,0)}, [_{(c,k,1)},]_{(c,k,1)}, \dots, [_{(c,k,k)},]_{(c,k,k)}$ is $k+1 < q$.

We introduce the following new brackets:

$$\begin{aligned}[c, 1,]_{c,1}, \dots, [c, q-1,]_{c,q-1}, \\ [_{(c,k,k+1)},]_{(c,k,k+1)}, \dots, [_{(c,k,q-1)},]_{(c,k,q-1)},\end{aligned}$$

for each $c \in \mathcal{Y}$ and $k < q - 1$. We now have an alphabet Γ_{qn} consisting of n groups of q pairs of brackets:

$$\begin{aligned} & \llbracket_c, \rrbracket_c, \llbracket_{c,1}, \rrbracket_{c,1}, \dots, \llbracket_{c,q-1}, \rrbracket_{c,q-1}, \\ & \llbracket_{(c,k,0)}, \rrbracket_{(c,k,0)}, \dots, \llbracket_{(c,k,q-1)}, \rrbracket_{(c,k,q-1)}, \end{aligned}$$

where $n = |\mathcal{Y}| \times (q + 1)$. Define a homomorphism $\psi: (\Gamma_{\tilde{\mathcal{Y}}_{q-1}}^*)^* \rightarrow \Gamma_{qn}^*$ by

$$\begin{aligned} \psi(\llbracket_c) &= \llbracket_c \llbracket_{c,1}, \\ \psi(\rrbracket_c) &= \rrbracket_{c,1} \rrbracket_{c,2} \dots \rrbracket_{c,q-1} \rrbracket_c, \\ \psi(\llbracket_{(c,k,0)}) &= \llbracket_{(c,k,0)}, \\ \psi(\rrbracket_{(c,k,0)}) &= \llbracket_{(c,k,k+1)} \rrbracket_{(c,k,k+1)} \dots \llbracket_{(c,k,q-1)} \rrbracket_{(c,k,q-1)} \rrbracket_{(c,k,0)}, \\ \psi(\llbracket_{(c,k,i)}) &= \llbracket_{(c,k,i)}, \\ \psi(\rrbracket_{(c,k,i)}) &= \rrbracket_{(c,k,i)} \quad \text{for } 1 \leq i \leq k. \end{aligned}$$

Now it is easy to see that ψ is injective and $\psi(R'')$ is a local subset of Γ_{qn}^* . Also, we can show that all $x \in (\Gamma_{\tilde{\mathcal{Y}}_{q-1}}^*)^*$ satisfy the following properties:

$$\begin{aligned} x \in D_{\tilde{\mathcal{Y}}_{q-1}} & \text{ if and only if } \psi(x) \in D_{qn}, \\ g(\psi(x)) & \rightsquigarrow^* \psi(\eta(x)). \end{aligned} \tag{20}$$

These properties combined ensure that

$$\eta(x) \in D_{\tilde{\mathcal{Y}}_{q-1}} \text{ if and only if } g(\psi(x)) \in D_{qn}. \tag{21}$$

Let $\chi: \Gamma_{qn}^* \rightarrow (\Gamma_{\tilde{\mathcal{Y}}_{q-1}}^*)^*$ be the alphabetic homomorphism that erases all new brackets. Clearly, χ restricted to the range of ψ is the inverse of ψ . Now observe

$$\chi(\psi(R'') \cap D_{qn} \cap g^{-1}(D_{qn})) = R'' \cap D_{\tilde{\mathcal{Y}}_{q-1}} \cap \eta^{-1}(D_{\tilde{\mathcal{Y}}_{q-1}}). \tag{22}$$

Indeed, if $x \in R''$, $\psi(x) \in D_{qn}$, and $g(\psi(x)) \in D_{qn}$, then $\chi(\psi(x)) = x \in R'' \cap D_{\tilde{\mathcal{Y}}_{q-1}}$ by (20), and $\eta(\chi(\psi(x))) = \eta(x) \in D_{\tilde{\mathcal{Y}}_{q-1}}$ by (21). Conversely, if $x \in R'' \cap D_{\tilde{\mathcal{Y}}_{q-1}}$ and $\eta(x) \in D_{\tilde{\mathcal{Y}}_{q-1}}$, then $\psi(x) \in \psi(R'') \cap D_{qn}$ by (20) and $g(\psi(x)) \in D_{qn}$ by (21), and so $x = \chi(\psi(x)) \in \chi(\psi(R'') \cap D_{qn} \cap g^{-1}(D_{qn}))$.

We obtain the statement of the lemma from (19) and (22) by taking $R = \psi(R'')$ and $h = h' \circ \chi$. \square

As before, in the direction (i) \Rightarrow (ii) of the above proof, $R \cap D_{qn} \cap g^{-1}(D_{qn}) = \psi(R'') \cap D_{qn} \cap g^{-1}(D_{qn})$ stands in one-one correspondence with the local set $K \subseteq \mathbb{T}_{\psi, \mathbf{x}}^3$. Thus, each derivation tree \mathbf{T} of a simple context-free tree grammar for M is uniquely represented by an element s of $R \cap D_{qn} \cap g^{-1}(D_{qn})$ such that $\mathbf{enc}(\mathbf{enc}_2(\mathbf{T}))$ is the image of s under a certain alphabetic homomorphism, and vice versa. This is exactly analogous to the situation with the original Chomsky-Schützenberger representation theorem.

As in the case of context-free languages, we can take a fixed Dyck language D_{2q} , instead of D_{qn} with varying n , and use a rational transduction to represent any string language that is the yield image of some $L \in \text{CFT}_{\text{sp}}(q-1)$:²⁷

Corollary 41. *For any $M \in \text{yCFT}_{\text{sp}}(q-1)$, there is a rational transduction τ such that $M = \tau(D_{2q} \cap g^{-1}(D_{2q}))$, where g is as defined in (16) with $n = 2$.*

10 Conclusion

In this paper, I have generalized Weir's [43] characterization of the string languages of tree-adjointing grammars to the string languages of simple context-free tree grammars of arbitrary fixed rank. I obtained this result via a natural generalization of the original Chomsky-Schützenberger theorem to simple context-free tree grammars. I represented derivation trees of simple context-free tree grammars as 3-dimensional trees, and proved this latter result as a general fact about simple context-free sets of m -dimensional trees, for arbitrary $m \geq 2$. This generality is of course an overkill for the purpose of obtaining my generalization of Weir's theorem, but it may be of independent interest. Moreover, all the complexity of the general case is essentially already present in the 3-dimensional case; proving only the special cases of the lemmas that are needed for the generalization of Weir's theorem will not be substantially simpler.

I emphasize that an important aspect of the original Chomsky-Schützenberger theorem is preserved both in my generalization of it to $\text{CFT}_{\text{sp}}(q-1)$ (the case $m = 2$ of Theorem 30) and in my generalization of Weir's theorem to $\text{yCFT}_{\text{sp}}(q-1)$ (Theorem 40): the representing set ($R \cap DT_{\mathcal{Y}}^2$ in the former and $R \cap D_{qn} \cap g^{-1}(D_{qn})$ in the latter) is in bijective correspondence with the set of derivation trees of the grammar. This holds with no restriction on the grammar except for the order of appearance of variables in the right-hand side of productions.²⁸

In order to define the m -dimensional yield of an $(m+1)$ -dimensional tree, I placed a restriction on the occurrences of the special label \mathbf{x} that serve as targets

²⁷ Also, when string languages over a fixed alphabet Σ with $|\Sigma| = k$ are considered, we can use $D_{q(k+2)}$ and a fixed alphabetic homomorphism h so that every $M \in \text{yCFT}_{\text{sp}}(q-1)$ can be written as $M = h(R \cap D_{q(k+2)} \cap g^{-1}(D_{q(k+2)}))$ for some *regular* set R . (This will require modification of the constructions used in several lemmas.) See [45] for an analogous characterization of string languages of q -MCFGs of rank r , and, e.g., [36] for the case of context-free languages.

²⁸ This contrasts with Arnold and Dauchet's [1] Chomsky-Schützenberger theorem for OI context-free tree languages, which relies on a normal form for OI context-free tree grammars that severely restricts the shape of the right-hand side of productions. I should also mention that their characterization of OI context-free tree languages does not seem to lead to a Weir-like characterization of indexed languages (i.e., OI macro languages), because in their version of Dyck tree languages, one occurrence of an opening bracket (corresponding to $(c, k, 0)$ in this paper) may be matched by an unbounded number of occurrences of closing brackets (corresponding to (c, k, i) with $i \geq 1$).

for substitution. If we wish to iterate the process of taking the yield, i.e., if we are interested in the yield of the yield of an $(m + 2)$ -dimensional tree, etc., we will need more than one variable as placeholders, with each variable providing targets for substitution at a different step of the iterative process of taking yields. Although I did not attempt to do so in this paper, it may be interesting to study the resulting hierarchy of classes of tree languages (and their yield images), the first three levels of the hierarchy being the local tree languages, the simple context-free tree languages, and the yields of the simple context-free sets of (well-labeled) 3-dimensional trees.²⁹ Since the function \mathbf{y}_m is an MSO-definable tree transduction (from $(m + 1)$ -ary trees to m -ary trees), we know that the tree languages at all levels of this hierarchy are within the tree-generating power of hyperedge replacement graph grammars, and the string languages that are their yield images are all multiple context-free languages (see, e.g., [10]).

References

1. Arnold, A., Dauchet, M.: Un theoreme de Chomsky-Schützenberger pour les forets algebriques. *CALCOLO* 14(2), 161–184 (1977)
2. Baldwin, W.A., Strawn, G.O.: Multidimensional trees. *Theoretical Computer Science* 84, 293–311 (1991)
3. Benoit, D., Demaine, E.D., Munro, I., Raman, R., Raman, V., Rao, S.S.: Representing trees of higher degree. *Algorithmica* 43, 275–292 (2005)
4. Chomsky, N., Schützeberger, M.P.: The algebraic theory of context-free languages. In: Braffort, P., Hirschberg, D. (eds.) *Computer Programming and Formal Systems*, pp. 118–161. North-Holland, Amsterdam (1963)
5. Chomsky, N.: Formal properties of grammars. In: Luce, R.D., Bush, R.R., Galanter, E. (eds.) *Handbook of Mathematical Psychology, Volume II*, pp. 323–418. John Wiley and Sons, New York (1963)
6. Comon, H., Dauchet, M., Gilleron, R., Jacquemard, F., Lugiez, D., Löding, C., Tison, S., Tommasi, M.: *Tree Automata Techniques and Applications*. Available on <http://tata.gforge.inria.fr/> (2008), release November 18, 2008
7. Eilenberg, S.: *Automata, Languages, and Machines*, vol. A. Academic Press, New York (1974)
8. Engelfreit, J., Maneth, S.: Tree languages generated by context-free graph grammars. In: Ehrig, H., Engels, G., Kreowski, H.J., Rozenberg, G. (eds.) *Theory and Application of Graph Transformations: 6th International Workshop, TAGT'98*. pp. 15–59. Springer, Berlin (2000)
9. Engelfriet, J., Schmidt, E.M.: IO and OI, part I. *The Journal of Computer and System Sciences* 15, 328–353 (1977)
10. Engelfriet, J.: Context-free graph grammars. In: Rozenberg, G., Salomaa, A. (eds.) *Handbook of Formal Languages, Volume 3: Beyond Words*, pp. 125–213. Springer, Berlin (1997)
11. Fischer, M.J.: *Grammars with Macro-Like Productions*. Ph.D. thesis, Harvard University (1968)
12. Fujiyoshi, A., Kasai, T.: Spinal-formed context-free tree grammars. *Theory of Computing Systems* 33, 59–83 (2000)

²⁹ Rogers [32] looked at a connection between multi-dimensional trees (with his notion of yield) and Weir’s [44] control language hierarchy.

13. Joshi, A.K., Schabes, Y.: Tree-adjoining grammars. In: Rozenberg, G., Salomaa, A. (eds.) *Handbook of Formal Languages*, vol. 3, pp. 69–123. Springer, Berlin (1997)
14. Kanazawa, M.: The convergence of well-nested mildly context-sensitive grammar formalisms (July 2009), an invited talk given at the 14th Conference on Formal Grammar, Bordeaux, France. Slides available at <http://research.nii.ac.jp/~kanazawa/>.
15. Kanazawa, M.: The pumping lemma for well-nested multiple context-free languages. In: Diekert, V., Nowotka, D. (eds.) *Developments in Language Theory: 13th International Conference, DLT 2009*. pp. 312–325. Springer, Berlin (2009)
16. Kanazawa, M.: Labeled bracketings and the Chomsky-Schützenberger theorem (May 2010), <http://makotokanazawa.blogspot.jp/2010/05/labeled-bracketings-and-chomsky.html>, blog post
17. Kanazawa, M.: Multi-dimensional trees and a Chomsky-Schützenberger-Weir representation theorem for simple context-free tree grammars. NII Technical Report NII-2013-003E, National Institute of Informatics (November 2013)
18. Kanazawa, M.: A generalization of linear indexed grammars equivalent to simple context-free tree grammars. In: *Formal Grammar, FG 2014*. Springer, Berlin (to appear)
19. Kanazawa, M., Salvati, S.: The copying power of well-nested multiple context-free grammars. In: Dediu, A.H., Fernau, H., Martín-Vide, C. (eds.) *Language and Automata Theory and Applications, Fourth International Conference, LATA 2010*. pp. 344–355. Springer, Berlin (2010)
20. Kasprzik, A.: Two Equivalent Regularizations for Tree Adjoining Grammars. Master’s thesis, University of Tübingen (2007)
21. Kasprzik, A.: Making finite-state methods applicable to languages beyond context-freeness via multi-dimensional trees. In: Piskorski, J., Watson, B.W., Yli-Jyrä, A. (eds.) *Finite-State Methods and Natural Language Processing*, pp. 98–109. IOS Press, Amsterdam (2009)
22. Kepser, S., Mönnich, U.: Closure properties of linear context-free tree languages with an application to optimality theory. *Theoretical Computer Science* 354(1), 82–97 (2006)
23. Kepser, S., Rogers, J.: The equivalence of tree adjoining grammars and monadic linear context-free tree grammars. *Journal of Logic, Language and Information* 20, 361–384 (2011)
24. Knuth, D.E.: *The Art of Computer Programming, Vol. I: Fundamental Algorithms*. Addison-Wesley, Reading, Mass., third edn. (1997)
25. Kozen, D.C.: *Automata and Computability*. Springer, New York (1997)
26. Libkin, L.: Logics for unranked trees: An overview. *Logical Methods in Computer Science* 2, 1–31 (2006)
27. Maletti, A., Engelfriet, J.: Strong lexicalization of tree adjoining grammars. In: *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics*. pp. 506–515. Association for Computational Linguistics (2012)
28. Matsubara, S., Kasai, T.: A characterization of TALs by the generalized Dyck language. *IEICE Transactions on Information and Systems (Japanese Edition)* J90-D, 1417–1427 (2007)
29. McNaughton, R., Papert, S.A.: *Counter-Free Automata*. MIT Press, Cambridge, Mass. (1971)
30. Mönnich, U.: Adjunction as substitution: An algebraic formulation of regular, context-free and tree adjoining languages (1997), arXiv:cmp-lg/9707012v1
31. Perrin, D.: Finite automata. In: van Leeuwen, J. (ed.) *Handbook of Theoretical Computer Science*, vol. B, pp. 1–57. Elsevier, Amsterdam (1990)

32. Rogers, J.: Syntactic structures as multi-dimensional trees. *Research on Language and Computation* 1, 265–305 (2003)
33. Rogers, J.: wMSO theories as grammar formalisms. *Theoretical Computer Science* 293, 291–320 (2003)
34. Rogers, J., Pullum, G.K.: Aural pattern recognition experiments and the subregular hierarchy. *Journal of Logic, Language and Information* 20, 329–342 (2011)
35. Rounds, W.: Mappings and grammars on trees. *Mathematical Systems Theory* 4, 257–287 (1970)
36. Salomaa, A.: *Formal Languages*. Academic Press, Orlando, Florida (1973)
37. Seki, H., Kato, Y.: On the generative power of multiple context-free grammars and macro grammars. *IEICE Transactions on Information and Systems* E91–D, 209–221 (2008)
38. Seki, H., Matsumura, T., Fujii, M., Kasami, T.: On multiple context-free grammars. *Theoretical Computer Science* 88(2), 191–229 (1991)
39. Sorokin, A.: Monoid automata for displacement context-free languages. In: *ESSLLI Student Session 2013 Preproceedings*. pp. 158–167 (2013)
40. Takahashi, M.: Generalizations of regular sets and their application to a study of context-free languages. *Information and Control* 27, 1–36 (1975)
41. Thatcher, J.W.: Characterizing derivation trees of context-free grammars through a generalization of finite automata theory. *Journal of Computer and System Sciences* 1, 317–322 (1967)
42. Thatcher, J.W.: Generalized² sequential machine maps. *Journal of Computer and System Sciences* 24, 339–367 (1970)
43. Weir, D.J.: *Characterizing Mildly Context-Sensitive Grammar Formalisms*. Ph.D. thesis, University of Pennsylvania (1988)
44. Weir, D.J.: A geometric hierarchy beyond context-free languages. *Theoretical Computer Science* 104(2), 235–261 (1992)
45. Yoshinaka, R., Kaji, Y., Seki, H.: Chomsky-Schützenberger-type characterization of multiple context-free languages. In: Dediú, A.H., Fernau, H., Martín-Vide, C. (eds.) *Language and Automata Theory and Applications, Fourth International Conference, LATA 2010*. pp. 596–607. Springer, Berlin (2010)