

# A Descriptive Characterization of Tree-Adjoining Languages (Project Note)

James Rogers

Dept. of Computer Science

Univ. of Central Florida, Orlando, FL, USA

## Abstract

Since the early Sixties and Seventies it has been known that the regular and context-free languages are characterized by definability in the monadic second-order theory of certain structures. More recently, these descriptive characterizations have been used to obtain complexity results for constraint- and principle-based theories of syntax and to provide a uniform model-theoretic framework for exploring the relationship between theories expressed in disparate formal terms. These results have been limited, to an extent, by the lack of descriptive characterizations of language classes beyond the context-free. Recently, we have shown that tree-adjoining languages (in a mildly generalized form) can be characterized by recognition by automata operating on three-dimensional tree manifolds, a three-dimensional analog of trees. In this paper, we exploit these automata-theoretic results to obtain a characterization of the tree-adjoining languages by definability in the monadic second-order theory of these three-dimensional tree manifolds. This not only opens the way to extending the tools of model-theoretic syntax to the level of TALs, but provides a highly flexible mechanism for defining TAGs in terms of logical constraints.

## 1 Introduction

In the early Sixties Büchi (1960) and Elgot (1961) established that a set of strings was regular iff it was definable in the weak monadic second-order theory of the natural numbers with successor ( $wS1S$ ). In the early Seventies an extension to the context-free languages was obtained by Thatcher and Wright (1968) and Doner (1970) who established that the CFLs were all and only the sets of strings forming the yield of sets of finite trees definable in the weak

monadic second-order theory of multiple successors ( $wSnS$ ). These descriptive characterizations have natural application to constraint- and principle-based theories of syntax. We have employed them in exploring the language-theoretic complexity of theories in GB (Rogers, 1994; Rogers, 1997b) and GPSG (Rogers, 1997a) and have used these model-theoretic interpretations as a uniform framework in which to compare these formalisms (Rogers, 1996). They have also provided a foundation for an approach to principle-based parsing via compilation into tree-automata (Morawietz and Cornell, 1997). Outside the realm of Computational Linguistics, these results have been employed in theorem proving with applications to program and hardware verification (Henriksen et al., 1995; Biehl et al., 1996; Kelb et al., 1997). The scope of each of these applications is limited, to some extent, by the fact that there are no such descriptive characterizations of classes of languages beyond the context-free. As a result, there has been considerable interest in extending the basic results (Mönnich, 1997; Volger, 1997) but, prior to the work reported here, the proposed extensions have not preserved the simplicity of the original results.

Recently, in (Rogers, 1997c), we introduced a class of labeled three-dimensional tree-like structures (three-dimensional tree manifolds—3-TM) which serve simultaneously as the derived and derivation structures of Tree Adjoining-Grammars (TAGs) in exactly the same way that labeled trees can serve as both derived and derivation structures for CFGs. We defined a class of automata over these structures that are a generalization of tree-automata (which are, in turn, an analogous generalization of ordinary finite-state automata over strings) and showed that the class of tree manifolds rec-

ognized by these automata are exactly the class of tree manifolds generated by TAGs if one relaxes the usual requirement that the labels of the root and foot of an auxiliary tree and the label of the node at which it adjoins all be identical.

Thus there are analogous classes of automata at the level of labeled three-dimensional tree manifolds, the level of labeled trees and at the level of strings (which can be understood as two- and one-dimensional tree manifolds) which recognize sets of structures that yield, respectively, the TALs, the CFLs, and the regular languages. Furthermore, the nature of the generalization between each level and the next is simple enough that many results lift directly from one level to the next. In particular, we get that the recognizable sets at each level are closed under union, intersection, relative complement, projection, cylindrification, and determinization and that emptiness of the recognizable sets is decidable. These are exactly the properties one needs to establish that recognizability by the automata over a class of structures characterizes satisfiability of monadic second-order formulae in the language appropriate for that class. Thus, just as the proofs of closure properties lift directly from one level to the next, Doner's and Thatcher and Wright's proofs that the recognizable sets of trees are characterized by definability in  $wSnS$  lift directly to a proof that the recognizable sets of three-dimensional tree manifolds are characterized by definability in their weak monadic second-order theory (which we will refer to as  $wSnT3$ ).

In this paper we carry out this program. In the next section we introduce 3-TMs, our uniform notion of automaton over tree manifolds of arbitrary (finite) dimension and indicate the nature of the dimension-independent proofs of closure properties. In Section 3 we introduce  $wSnT3$ , the weak monadic second-order theory of  $n$ -branching 3-TM, and sketch the proof that the sets definable in  $wSnT3$  are exactly those recognizable by 3-TM automata. This, when coupled with the characterization of TALs in Rogers (1997c), gives us our descriptive characterization of TALs: a set of strings is generated by a TAG (modulo the generalization of Rogers (1997c)) iff it is the (string) yield of a set of 3-TM definable in  $wSnT3$ . Finally, in Sec-

tion 4 we look at how working in  $wSnT3$  allows a potentially more transparent means of defining TALs and, in particular, a simplified treatment of constraints on modifiers in TAGs. Due to the limited length of this note, many of the details are omitted. The reader is directed to (Rogers, 1998) for a more complete treatment.

## 2 Tree Manifolds and Automata

Tree manifolds are a generalization to arbitrary dimensions of Gorn's *tree domains* (Gorn, 1967). A tree domain is a set of node address drawn from  $\mathbb{N}^*$  (that is, a set of strings of natural numbers) in which  $\varepsilon$  is the address of the root and the children of a node at address  $w$  occur at addresses  $w0, w1, \dots$ , in left-to-right order. To be well formed, a tree domain must be downward closed wrt to domination, which corresponds to being prefix closed, and left sibling closed in the sense that if  $wi$  occurs then so does  $wj$  for all  $j < i$ . In generalizing these, we can define a one-dimensional analog as *string domains*: downward closed sets of natural numbers interpreted as string addresses. From this point of view, the address of a node in a tree domain can be understood as the sequence of string addresses one follows in tracing the path from the root to that node. If we represent  $\mathbb{N}$  in unary (with  $n$  represented as  $1^n$ ) then the downward closure property of string domains becomes a form of prefix closure analogous to downward closure wrt domination in tree domains, tree domains become sequences of sequences of '1's, and the left-closure property of tree domains becomes a prefix closure property for the embedded sequences.

Raising this to higher dimensions, we obtain, next, a class of structures in which each node expands into a (possibly empty) tree. A *three-dimensional tree manifold* (3-TM), then, is set of sequences of tree addresses (that is, addresses of nodes in tree domains) tracing the paths from the root of one of these structures to each of the nodes in it. Again this must be downward closed wrt domination in the third dimension, equivalently wrt prefix, the sets of tree addresses labeling the children of any node must be downward closed wrt domination in the second dimension (again wrt to prefix), and the sets of string addresses labeling the children of any node in any of these trees must be downward

closed wrt domination in the first dimension (left-of, and, yet again, prefix). Thus 3-TM, tree domains (2-TM), and string domains (1-TM) can be defined uniformly as  $d^{\text{th}}$ -order sequences of '1's which are hereditarily prefix closed. We will denote the set of all 3-TM as  $\mathbb{T}^d$ . For any alphabet  $\Sigma$ , a  $\Sigma$ -labeled  $d$ -dimensional tree manifold is a pair  $\langle T, \tau \rangle$  where  $T$  is a  $d$ -TM and  $\tau : T \rightarrow \Sigma$  is an assignment of labels in  $\Sigma$  to the nodes in  $T$ . We will denote the set of all  $\Sigma$ -labeled  $d$ -TM as  $\mathbb{T}_\Sigma^d$ .

Mimicking the development of tree manifolds, we can define automata over labeled 3-TM as a generalization of automata over labeled tree domains which, in turn, can be understood as an analogous generalization of ordinary finite-state automata over strings (labeled string domains). A  $d$ -TM automaton with state set  $Q$  and alphabet  $\Sigma$  is a finite set:

$$\mathcal{A}^d \subseteq \Sigma \times Q \times \mathbb{T}_Q^{d-1}.$$

The interpretation of a tuple  $\langle \sigma, q, \mathcal{T} \rangle \in \mathcal{A}^d$  is that if a node of a  $d$ -TM is labeled  $\sigma$  and  $\mathcal{T}$  encodes the assignment of states to its children, then that node may be assigned state  $q$ . A run of an  $d$ -TM automaton  $\mathcal{A}$  on a  $\Sigma$ -labeled  $d$ -TM  $\mathcal{T} = \langle T, \tau \rangle$  is an assignment  $r : T \rightarrow Q$  of states in  $Q$  to nodes in  $T$  in which each assignment is licensed by  $\mathcal{A}$ . If we let  $Q_0 \subseteq Q$  be any set of *accepting states*, then the set of (finite)  $\Sigma$ -labeled  $d$ -TM recognized by  $\mathcal{A}$ , relative to  $Q_0$ , is that set for which there is a run of  $\mathcal{A}$  that assigns the root a state in  $Q_0$ . A set of  $d$ -TM is *recognizable* iff it is  $\mathcal{A}(Q_0)$  for some  $d$ -TM automaton  $\mathcal{A}$  and set of accepting states  $Q_0$ .

The strength of the uniform definition of  $d$ -TM automata is that many, even most, properties of the sets they recognize can be proved uniformly—independently of their dimension. It is easy to see that in the typical “cross-product” construction of the proof of closure under intersection, for instance, the dimensionality of the TMs is a parameter that determines the type of the objects being manipulated but does not affect the manner of their manipulation. Uniform proofs can be obtained for closure of recognizable sets under determinization (in a bottom-up sense), projection, cylindrification, Boolean operations and for decidability of emptiness.

### 3 wSnT3

We are now in a position to build relational structures on  $d$ -dimensional tree manifolds. Let  $T_n^d$  be the *complete  $n$ -branching  $d$ -TM*—that in which every point has a child structure that has depth  $n$  in all its  $(d - 1)$  dimensions. Let

$$\mathbb{T}_n^3 \stackrel{\text{def}}{=} \langle T_n^3, \triangleleft_1, \triangleleft_2, \triangleleft_3 \rangle$$

where, for all  $x, y \in T_n^3$ ,  $x \triangleleft_i y$  iff  $x$  is the immediate predecessor of  $y$  in the  $i^{\text{th}}$ -dimension.

The *weak monadic second-order language* of  $\mathbb{T}_n^3$  includes constants for each of the relations (we let them stand for themselves), the usual logical connectives, quantifiers and grouping symbols, and two countably infinite sets of variables, one ranging over individuals (for which we employ lowercase) and one ranging over finite subsets (for which we employ uppercase). If  $\varphi(x_1, \dots, x_n, X_1, \dots, X_m)$  is a formula of this language with free variables among the  $x_i$  and  $X_j$ , then we will assert that it is satisfied in  $\mathbb{T}_n^3$  by an assignment  $\mathbf{s}$  (mapping the ' $x_i$ 's to individuals and ' $X_j$ 's to finite subsets) with the notation  $\mathbb{T}_n^3 \models \varphi[\mathbf{s}]$ . The set of all sentences of this language that are satisfied by  $\mathbb{T}_n^3$  is the *weak monadic second-order theory* of  $\mathbb{T}_n^3$ , denoted wSnT3.

A set  $\mathbb{T}$  of  $\Sigma$ -labeled 3-TM is definable in wSnT3 iff there is a formula  $\varphi_{\mathbb{T}}(X_T, X_\sigma)_{\sigma \in \Sigma}$ , with free variables among  $X_T$  (interpreted as the domain of a tree) and  $X_\sigma$  for each  $\sigma \in \Sigma$  (interpreted as the set of  $\sigma$ -labeled points in  $T$ ), such that

$$\langle T, \tau \rangle \in \mathbb{T} \iff \mathbb{T}_n^3 \models \varphi_{\mathbb{T}} [X_i \mapsto T, X_\sigma \mapsto \{p \mid \tau(p) = \sigma\}].$$

It should be reasonably easy to see that any recognizable set can be defined by encoding the local TM of an accepting automaton in formulae in which the labels and states occur as free variables and then requiring every node to satisfy one of those formulae. One then requires the root to be labeled with an accepting state and “hides” the states by existentially binding them.

The proof that every set of trees definable in wSnT3 is recognizable, while a little more involved, is just a lift of the proofs of Doner and Thatcher and Wright. The initial step is to show that every formula in the language of wSnT3

can be reduced to equivalent formulae in which only set variables occur and which employ only the predicates  $X \subseteq Y$  (with the obvious interpretation) and  $X \triangleleft_i Y$  (satisfied iff  $X$  and  $Y$  are both singleton and the sole element of  $X$  stands in the appropriate relation to the sole element of  $Y$ ). It is easy to construct 3-TM automata (over the alphabet  $\mathcal{P}(\{X, Y\})$ , where  $\mathcal{P}$  denotes power set) which accept trees encoding satisfying assignments for these atomic formulae. The extension to arbitrary formulae (over these atomic formulae) can then be carried out by induction on the structure of the formulae using the closure properties of the recognizable sets.

#### 4 Defining TALs in wSnT3

The signature of wSnT3 is inconvenient for expressing linguistic constraints. In particular, one of the strengths of the model-theoretic approach is the ability to define long-distance relationships without having to explicitly encode them in the labels of the intervening nodes. We can extend the immediate predecessor relations to relations corresponding to (proper) *above* (within the 3-TM), *domination* (within a tree), and *precedence* (within a set of siblings) using:

$$x \bar{\triangleleft}_i y \stackrel{\text{def}}{\iff} x \neq y \wedge (\exists X)[X(x) \wedge X(y) \wedge (\forall z)[X(z) \rightarrow (z \approx y \vee (\exists! z')[X(z') \wedge z \triangleleft_i z'])]]$$

Which simply asserts that there is a sequence of (at least two) points linearly ordered by  $\triangleleft_i$  in which  $x$  precedes  $y$ .

To extend these through the entire structure we have to address the fact that the two dimensional yield of a 3-TM is not well defined—there is nothing that determines which leaf of the tree expanding a node dominates the subtree rooted at that node. To resolve this, we extend our structures to include a set  $H$  picking out exactly one head in each set of siblings, with the “foot” of a tree being that leaf reached from the root by a path of all heads. Given  $H$ , it is possible to define  $\triangleleft_2^+$  and  $\triangleleft_1^+$ , variations of dominance and precedence<sup>1</sup> that are inherited by substructures in the appropriate way. At the same time, it is convenient to include the labels explicitly in the structures. A headed  $\Sigma$ -labeled 3-TM, then, is

<sup>1</sup>Of course  $\triangleleft_3^+$  is just  $\bar{\triangleleft}_3$ .

a structure:

$$\langle T, \triangleleft_i, \bar{\triangleleft}_i, \triangleleft_i^+, H, P_\sigma \rangle_{1 \leq i \leq 3, \sigma \in \Sigma},$$

where  $T$  is a rooted, connected subset of  $T_n^3$  for some  $n$ .

With this signature it is easy to define the set of 3-TM that captures a TAG in the sense that their 2-dimensional yields—the set of maximal points wrt  $\triangleleft_3^+$ , ordered by  $\triangleleft_2^+$  and  $\triangleleft_1^+$ —form the set of trees derived by the TAG. Note that obligatory (OA) and null (NA) adjoining constraints translate to a requirement that a node be (non-)maximal wrt  $\triangleleft_3^+$ . In our automata-theoretic interpretation of TAGs selective adjoining (SA) constraints are encoded in the states. Here we can express them directly: a constraint specifying the modifier trees which may adjoin to an N node, for instance, can be stated as a condition on the label of the root node of trees immediately below N nodes.

In general, of course, SA constraints depend not only on the attributes (the label) of a node, but also on the elementary tree in which it occurs and its position in that tree. Both of these conditions are actually expressions of the local context of the node. Here, again, we can express such conditions directly—in terms of the relevant elements of the node’s neighborhood. At least in some cases this seems likely to allow for a more general expression of the constraints, abstracting away from the irrelevant details of the context.

Finally, there are circumstances in which the primitive locality of SA constraints in TAGs is inconvenient. Schabes and Shieber (1994), for instance, suggest allowing multiple adjunctions of modifier trees to the same node on the grounds that selectional constraints hold between the modified node and each of its modifiers but, if only a single adjunction may occur at the modified node, only the first tree that is adjoined will actually be local to that node. They point out that, while it is possible to pass these constraints through the tree by encoding them in the labels of the intervening nodes, such a solution can have wide ranging effects on the overall grammar. As we noted above, the expression of such non-local constraints is one of the strengths of the model-theoretic approach. We can state them in a purely natural way—as a simple restriction on the types of the modifier

trees which can occur below (in the  $\leq_3^+$  sense) the modified node.

## 5 Conclusion

We have obtained a descriptive characterization of the TALs via a generalization of existing characterizations of the CFLs and regular languages. These results extend the scope of the model-theoretic tools for obtaining language-theoretic complexity results for constraint- and principle-based theories of syntax to the TALs and, carrying the generalization to arbitrary dimensions, should extend it to cover a wide range of mildly context-sensitive language classes. Moreover, the generalization is natural enough that the results it provides should easily integrate with existing results employing the model-theoretic framework to illuminate relationships between theories. Finally, we believe that this characterization provides an approach to defining TALs in a highly flexible and theoretically natural way.

## References

- M. Biehl, N. Klarlund, and T. Rauhe. 1996. Algorithms for guided tree automata. In *WIA '96*, LNCS 1260, London, Ontario.
- J. R. Büchi. 1960. Weak second-order arithmetic and finite automata. *Zeitschrift für mathematische Logik und Grundlagen der Mathematik*, 6:66–92.
- John Doner. 1970. Tree acceptors and some of their applications. *Journal of Computer and System Sciences*, 4:406–451.
- Calvin C. Elgot. 1961. Decision problems of finite automata design and related arithmetics. *Transactions of the American Mathematical Society*, 98:21–51.
- Saul Gorn. 1967. Explicit definitions and linguistic dominoes. In *Systems and Computer Science, Proceedings of the Conference held at Univ. of Western Ontario, 1965*. Univ. of Toronto Press.
- J. G. Henriksen, J. Jensen, M. Jørgensen, N. Klarlund, R. Paige, T. Rauhe, and A. Sandhol. 1995. MONA: Monadic second-order logic in practice. In *TACAS '95*, LNCS 1019, Aarhus, Denmark.
- P. Kelb, T. Margaria, M. Mendler, and C. Gsotberger. 1997. MOSEL: A flexible toolset for monadic second-order logic. In *TACAS '97*, LNCS 1217, Enschede, The Netherlands.
- Uwe Mönnich. 1997. Adjunction as substitution: An algebraic formulation of regular, context-free and tree adjoining languages. In *Formal Grammar*, Aix-en-Provence, Fr.
- Frank Morawietz and Tom Cornell. 1997. Representing constraints with automata. In *Proceedings of the 35th Annual Meeting of the ACL*, pages 468–475, Madrid, Spain.
- James Rogers. 1994. *Studies in the Logic of Trees with Applications to Grammar Formalisms*. Ph.D. thesis, Department of Computer and Information Sciences, University of Delaware.
- James Rogers. 1996. A model-theoretic framework for theories of syntax. In *Proceedings of the 34th Annual Meeting of the ACL*, pages 10–16, Santa Cruz, CA.
- James Rogers. 1997a. “Grammarless” phrase structure grammar. *Linguistics and Philosophy*, 20:721–746.
- James Rogers. 1997b. On descriptive complexity, language complexity, and GB. In *Specifying Syntactic Structures*, pages 157–184. CSLI Publications.
- James Rogers. 1997c. A unified notion of derived and derivation structures in TAG. In *Proceedings of the Fifth Meeting on Mathematics of Language MOL5 '97*, Saarbrücken, FRG.
- James Rogers. 1998. A descriptive characterization of tree-adjoining languages. Technical Report CS-TR-98-01, Univ. of Central Florida. Also available from the CMP-LG repository as paper number cmp-lg/9805008.
- Yves Schabes and Stuart M. Shieber. 1994. An alternative conception of tree-adjoining derivation. *Computational Linguistics*, 20(1):91–124.
- J. W. Thatcher and J. B. Wright. 1968. Generalized finite automata theory with an application to a decision problem of second-order logic. *Mathematical Systems Theory*, 2(1):57–81.
- Hugo Volger. 1997. Principle languages and principle based parsing. Technical Report 82, SFB 340, Univ. of Tübingen.