

# Model Checking Secondary Relations

Stephan Kepser, Uwe Mönnich, and Frank Morawietz

Seminar für Sprachwissenschaft,  
Universität Tübingen, Germany  
{kepser, um, frank}@sfs.uni-tuebingen.de

## 1 Introduction

Regarding linguistic data structures as relational structures makes them amenable to the techniques of model checking. The basic question in this area concerns the problem of how to devise efficient procedures that tell structures exhibiting a certain property from those that lack this property. As these properties are expressed by means of logical formulae, one can also regard the problem of model checking as a form of querying relational structures. Of special interest in this connection are formulae in the language of monadic second-order logic (MSO).

In the present paper we try to take advantage of a powerful generalisation of the classical result that the intersection of a context-free language and a regular one is a context-free language. The generalisation consists in defining a family of structures as context-free if it is a component of the least fixed-point of a system of equations over a finite set  $\mathcal{F}$  of operations. In the particular case of context-free languages the equations are expressed with the operations of union and concatenation. Using instead a distinguished set of projection and composition operations it becomes possible to characterise structures by means of appropriate systems of equations that are located on higher levels of the Chomsky hierarchy. Based on the particular set of operation symbols just mentioned the whole family of indexed languages can be accommodated within this framework.

The extension of the classical result concerning the intersection of context-free and regular languages depends on an important property the set  $T(\mathcal{F})$  of trees over  $\mathcal{F}$  and the associated evaluation  $val_{\mathcal{F}}$  which sends these trees into the intended structures has to satisfy. This property has been called MSO compatibility by Courcelle and Walukiewicz (1998) and requires of a (partial) mapping  $f$  from structures  $S(R)$  over the signature  $R$  into structures  $S(R')$  over the signature  $R'$  that for every MSO-

sentence  $\varphi$  one can produce a backwards translation  $f^{\sharp}(\varphi)$  such that

$$S \models f^{\sharp}(\varphi) \text{ iff } f(S) \models \varphi$$

for every structure  $S$  in the domain of  $f$ . Given the well-known fact that MSO definability on trees is equivalent to recognisability by finite tree automata the generalisation of the classical result follows immediately once the set of operators  $\mathcal{F}$  under consideration is MSO compatible.

In a series of papers (Kolb et al., 2000a; Kolb et al., 2000b; Michaelis et al., 2001; Morawietz and Mönnich, 2001) Mönnich, Morawietz, Kolb, and Michaelis have shown that the evaluation that interprets trees over projection and composition within the domain of structures familiar to linguists can be expressed as a simple form of MSO transduction ( $def_{\Delta}(MSO)$ ). Such a transduction defines the intended structure, i.e.  $val(t)$ , for  $t \in T(\mathcal{F})$ , within the input structure  $t$  on the basis of a finite set of MSO formulae written in the signature of the input structure. These defining formulae can then be used for the backwards translation ( $def_{\Delta}(MSO)^{-1}$ ) of a property that is expressed by an MSO formula over the target signature. Succinctly: MSO transductions are MSO compatible.

As it turns out, this relationship between the hierarchical structure of trees over composition and projection and their intended interpretation provides the foundation for a very flexible model checking procedure. Suppose one is dealing with a class of structures that exhibit a certain number of secondary relations. For concreteness assume that these relations indicate the sort of context-sensitive dependencies which have been at the focus of attention of linguists. As has been noticed since the beginning of formal language theory certain grammatical phenomena like morphological congruences (e.g., in Bambara) and cross-serial dependencies between case markings (e.g., in Swiss German) are outside the realm of context-free languages and need

for their descriptive analysis a (limited) amount of contextual information (Shieber, 1985). Due to this character of context-sensitivity that comes with a range of grammatical constructions, even MSO logic, normally considered a powerful query language, is too weak to capture these phenomena.

Taking our inspiration from the concept of MSO compatibility we regard grammatical categories as basic constants of a many-sorted algebra with a distinguished set of composition and projection symbols. Through the explicit introduction of these operations it becomes possible to turn the data models of contemporary syntactic theories into a kind of labeled tree structures that can either be generated by regular tree grammars or are identifiable with collections of finite trees specifiable by formulae of MSO logic.

In the particular case of the verbal complex of Swiss German mentioned above it is easy to describe in MSO terms the two verbal and nominal clusters, respectively. What is problematic from the point of view of regularity is the set of fixed syntactic and semantic relations between the verbal elements and their overtly case-marked arguments. In other words, an MSO specification of these bipartite structures would return – regarding the MSO specification as a yes/no query – structures that do not satisfy the particular set of cross-serial dependencies characteristic of the instance of context-sensitivity under discussion. Despite this lack of expressive power of the chosen query language the general approach adumbrated above is remarkably effective in filtering out the syntactic “noise” from the query result. Since the explicit algebraic structures are elements of a regular family of trees it is again easy to produce an MSO formula that characterises exactly these cross-serial dependencies among the explicit structures that were out of the reach of the query language on the intended linguistic level.

Generalizing from the particular problem of cross-serial dependencies in natural languages the impact of the classical result from formal language theory mentioned above can be described as follows. If a set of operations  $\mathcal{F}$  and the associated interpretation function  $val_{\mathcal{F}}$  satisfy the condition of MSO compatibility, the subset of structures within a context-free language  $\mathcal{L}$  (in the general sense) that fulfill a certain MSO formula  $\varphi$  is context-free. In symbols:

$$\{S \mid S \models \varphi \wedge S \in \mathcal{L}\} \in CF$$

Relying again on the fact that MSO definability is equivalent to recognizability by finite tree automata the subset of  $\mathcal{L}$  specified by the formula  $\varphi$  can be given an efficient regular description on the level of trees  $T(\mathcal{F})$ .

It has been shown that this method of regularising queries of context-sensitive structures can be adapted to all grammatical phenomena that fall within the reach of current linguistic theories (see the list of papers cited above). Since our data model is firmly entrenched in the linguistic tradition where trees with a limited amount of cross-serial dependencies play a prominent role, we are able to restrict our attention to two constructors denoting the familiar operations of composition and projection. This advantage which is provided by the considerable reduction of the set of primitive constructors does not lead directly to a family of canonical expressions that suits our purposes. As was noted above a query whose expressive power does not go beyond MSO is too weak to specify an answer set displaying the sort of dependencies so characteristic of natural language structure. It is therefore necessary to translate the result of the query into a family of trees that can be checked by a suitable constraint formula for the intended dependencies. We will show below that this translation of the first step is tightly controlled by the constraint formula. Using the constraint formula as a template for the translation process allows us to avoid the problem of context-sensitive parsing without being forced to consider the unbounded set of “lifted” expressions denoting the same tree.

The idea of using composition and projection as operations on trees is a special case of a general approach developed by Mezei and Wright (1967) in which regular tree languages denote subsets of arbitrary algebras. Of particular relevance for the present application to context-sensitive query problems have been the contributions of Courcelle (1990) to the interaction between graph operations and MSO. Courcelle has devised a primitive set of operations such that any finite graph can be considered as the value of a term that is constructed from (symbols for) these primitive operations. In a recent paper Courcelle and Knapik (2002) prove that the mapping which associates a term  $t$  over a complete set of graph operations with its evaluation  $val(t)$  is an MSO-transduction (Proposition 2.5).

The method of turning the classical result from formal language theory into a powerful model checking procedure can be put to use in the con-

text of recent attempts to specify a common logical level for linguistic databases. As has been emphasised by Cotton and Bird (2002) the proliferation of linguistic databases with their bewildering diversity of formats and software tools makes it necessary to integrate them into a general multilayer annotation system. For the special case of treebank formats the authors show how they can be mapped onto the annotation graph model serving as a common logical level. Since the annotation graph model can be regarded as a special type of relational structure it is again easy to verify that the mapping from the entries of a treebank into annotation graphs is an interpretation along this lines of an MSO-transduction. The inverse mapping from the annotation graph structures into the treebank format produces tree-like graphs with crossing lines if the annotation graphs contain equivalence classes or cross references of edges. The method of the present paper of how to exploit the universal algebraic version of the closure property of context free languages with respect to the intersection with regular languages is then applicable to the resulting trees with secondary relations and we are thus able to "lift" the problem of model checking from the level of annotation graphs to the level of trees without crossing lines.

## 2 Hierarchical Decomposition and Model Checking

The aim of the paper is to provide a way to check secondary relations in a context-free database. A database in our sense is just a set of relational structures. As mentioned above, we use monadic second-order logic as query language. A query is therefore an MSO-sentence, and the answer to a query is the set of all those structures in the database for which this formula is true. But since MSO is restricted to context-free phenomena, we need a device to specify the (mildly) context-sensitive secondary relations a linguist may be interested in. Obviously, the linguist has to specify a grammar that generates the structures he is interested in. This requires the grammar formalism to be more expressive than context-free (string) grammars. The largest class of grammars suitable for our approach are linear context-free tree grammars.

**Definition 1** [Context-Free Tree Grammar] Let  $\mathcal{S}$  be a singleton set of sorts. Then a *context-free tree grammar* (CFTG) for  $\mathcal{S}$  is a 5-tuple  $\Gamma =$

$\langle \Sigma, F, S, X, P \rangle$ , where  $\Sigma$  and  $F$  are ranked alphabets of *inoperatives* and *operatives* over  $\mathcal{S}$ , respectively.  $S \in F$  is the start symbol,  $X$  is a countable set of variables, and  $P$  is a set of productions. Each  $p \in P$  is of the form  $F(x_1, \dots, x_n) \longrightarrow t$  for some  $n \in \mathbb{N}$ , where  $F \in F_n$ ,  $x_1, \dots, x_n \in X$ , and  $t \in T(\Sigma \cup F, \{x_1, \dots, x_n\})$ .<sup>1</sup> The grammar is *linear*, iff each variable occurs at most once in the left hand side and at most once in the right hand side of a rule.

Intuitively, an application of a rule of the form  $F(x_1, \dots, x_n) \longrightarrow t$  "rewrites" a tree rooted in  $F$  as the tree  $t$  with its respective variables substituted by  $F$ 's daughters. A context-free tree grammar generates elements of a tree substitution algebra  $DT(\Sigma, X)$ .

A CFTG  $\Gamma = \langle \Sigma, F, S, X, P \rangle$  with  $F_n = \emptyset$  for  $n \neq 0$  is called a *regular tree grammar* (RTG). Since RTGs always just substitute some tree for a leaf-node, it is easy to see that they can only generate recognizable sets of trees, *a fortiori* context-free string languages (Mezei and Wright, 1967). If  $F_n$  is non-empty for some  $n \neq 0$ , that is, if we allow the *operatives* to be parameterized by variables, however, the situation changes. CFTGs in general are capable of generating sets of structures, the *yields* of which belong to the subclass of context-sensitive languages known as the *indexed* languages.

In order to be able to find the desired context-sensitive relations, it is necessary that the grammar used is such that it generates all those trees which embody non-context-free relations. It need not be a grammar for a single query. But it should usually not be a general grammar for the whole database either, because it will be used as a filter.

In order to apply the general approach adumbrated in the introduction to the kind of secondary relations that can be accommodated within the framework of context-free tree grammars we need an appropriate set of operations that subtend the necessary hierarchical decomposition. The intuition here is that the basic assumptions about the operations of a tree grammar, namely tree substitution and argument insertion, are made explicit. We make them visible by inserting the "control" information which allows us to code the resulting structures with regular means, i.e., regular tree grammars or finite-state tree automata and therefore with MSO logic. The intuition behind the LIFTing process is that each term compactly encodes information such as com-

<sup>1</sup> $T(\Sigma, X)$  stands for the set of trees (or terms) over the ranked alphabet  $\Sigma$  and the variables in  $X$  built in the usual way.

position and concatenation.

In the following, we will briefly describe LIFTing on an informal level. All technical details, in particular concerning many-sorted signatures, can be found in a paper by Mönnich (1999). Any *context-free* tree grammar  $\Gamma$  for a singleton set of sorts  $\mathcal{S}$  can be transformed into a *regular* tree grammar  $\Gamma^L$  for the set of sorts  $\mathcal{S}^*$ , which characterizes a (necessarily recognizable) set of trees encoding the instructions necessary to convert them by means of a unique homomorphism  $h$  into the ones the original grammar generates (Maibaum, 1974). This unique homomorphism  $h$  is nothing else but the evaluation mapping *val* discussed above. The “LIFTing” is achieved by constructing for a given single-sorted signature  $\Sigma$  a new, derived alphabet (an  $\mathcal{S}^*$ -sorted signature)  $\Sigma^L$ , and by translating the terms over the original signature into terms of the derived one via a primitive recursive procedure. The LIFT-operation takes a term in  $T(\Sigma, X_k)$  and transforms it into one in  $T(\Sigma^L, k)$ .<sup>2</sup> Intuitively, the LIFTing eliminates variables and composes functions with their arguments explicitly, e.g., a term  $f(a, b) = f(x_1, x_2) \circ (a, b)$  is lifted to the term  $c(c(f, \pi_1, \pi_2), a, b)$ . The old function symbol  $f$  now becomes a constant, the variables are replaced with appropriate projection symbols and the only remaining non-nullary alphabet symbols are the explicit composition symbols  $c$ . The trees over the derived “LIFTed” signature consisting of the old linguistic symbols together with the new projection and composition symbols form the carrier of a free tree algebra  $T_D$ .

Our main result provides a basis for a definition of the linguistically meaningful structures of the tree substitution algebra within the trees of the LIFTed algebra. Actually, it consists of a variant of the classical technique of interpreting one relational structure inside another one. The particular variant we use is due to Courcelle (1997) and interprets the domain and the relations on the substitution trees by means of suitable MSO formulae written in the signature of the derived algebra.

**Proposition 2** *The evaluation from the free derived algebra  $T_D$  into the tree substitution algebra  $DT(\Sigma, X)$  is an MSO-transduction.*

The idea for the proof of this proposition is due to Kolb (1999). He was able to analyse the elements

<sup>2</sup>Since  $\mathcal{S}$  is a singleton set of sorts, we can identify  $\mathcal{S}^*$  with  $\mathbb{N}$ . By  $T(\Sigma^L, k)$  we denote the set of all trees over  $\Sigma^L$  which are of sort  $k$ .

of  $T_D$  in such a way that the mapping from the free derived tree algebra into the tree substitution algebra can be simulated by a tree walking automaton. The walks on a tree this automaton accepts connect nodes that satisfy the relations on the substitution trees. As explained in detail in the paper by Morawietz and Mönnich (2001), these walks are specifiable by suitable MSO formulae, thereby providing the desired logical definitions of the target relations. As an immediate consequence one has the following corollary.

**Corollary 3** *The transformation of trees in  $T_D$  by means of composing and projecting subterms is MSO-compatible.*

The backwards translation  $def_\Delta(MSO)^{-1}$ , the inverse of the MSO-transduction, provides a way to “lift” the input query: The relation symbols of the input signature can be replaced by their images under  $def_\Delta(MSO)^{-1}$ . As a result we receive an MSO-query in the “lifted” signature.

According to the corollary any MSO query  $\varphi$  addressed at the database can be translated into a query  $val^\sharp(\varphi)$ , i.e., the result of replacing the relation symbols occurring in  $\varphi$  by their images under  $def_\Delta(MSO)^{-1}$ , phrased by means of the derived vocabulary incorporating the composition and the projection symbols. By the well-known equivalence between tree automata and MSO formulae  $val^\sharp(\varphi)$  has a translation into a corresponding tree automaton. The same transformation applied to the context-free tree grammar in the background of the supposed data base produces another tree automaton. Intersection of these two automata produces an automaton that accepts only structures that are in conformity with both the background grammar and the “lifted” query  $val^\sharp(\varphi)$ . The intended answer set to the original query contains those elements in the language of the final automaton whose evaluations via  $def_\Delta(MSO)$  pass a simple membership test against the given data base. Figure 1 gives an overview of our approach as we described it above. There also exists a prototypical implementation which interleaves the necessary processes. Due to space limitations, we cannot give any details in this extended abstract.

In this paper, we presented an approach to checking context-sensitive relations with purely regular means. At the heart of this approach lies the insight that lifting a context-free tree grammar results in a regular tree grammar, which, since it is regular, can

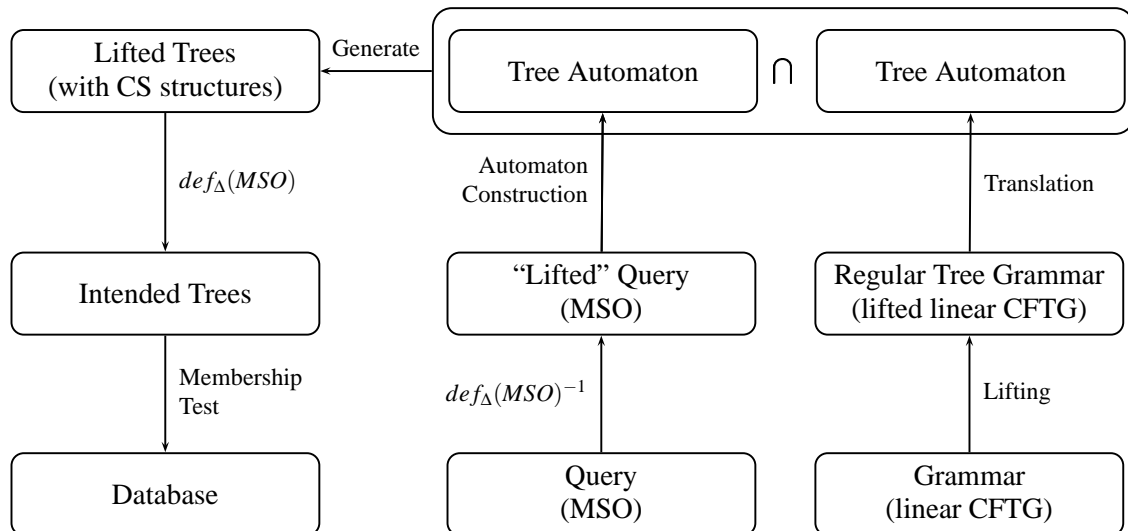


Figure 1: Overview of the approach

again be handled by monadic second order logic and its associated automata theory. The seeming contradiction of using regular means to query mildly context-sensitive relations gets resolved by the old result (see, e.g., Courcelle (1990)) that the application of MSO-definable transductions on MSO-definable structures results in structures that may no longer be MSO-expressible.

## References

- Scott Cotton and Steven Bird. 2002. An integrated framework for treebanks and multilayer annotations. In *Proceedings LREC 2002*.
- Bruno Courcelle and Teodor Knapik. 2002. The evaluation of first-order substitution is monadic second-order compatible. To appear in: *TCS*.
- Bruno Courcelle and Igor Walukiewicz. 1998. Monadic second-order logic, graph coverings and unfoldings of transition systems. *Annals of Pure and Applied Logic*, 92:35–62.
- Bruno Courcelle. 1990. Graph rewriting, an algebraic and logic approach. In Jan van Leeuwen, editor, *Handbook of Theoretical Computer Science*, volume B, pages 193–242. Elsevier.
- Bruno Courcelle. 1997. The expression of graph properties and graph transformations in monadic second-order logic. In G. Rozenberg, editor, *Handbook of Graph Grammars and Computing by Graph Transformation. Vol. I: Foundations*, pages 313–400. World Scientific.
- Hans-Peter Kolb and Uwe Mönnich, editors. 1999. *The Mathematics of Syntactic Structure*. Mouton de Gruyter.
- Hans-Peter Kolb, Jens Michaelis, Uwe Mönnich, and Frank Morawietz. 2000a. An operational and denotational approach to non-context-freeness. To appear in: *TCS*.
- Hans-Peter Kolb, Uwe Mönnich, and Frank Morawietz. 2000b. Descriptions of cross-serial dependencies. *Grammars*, 3(2/3):189–216.
- Hans-Peter Kolb. 1999. Macros for minimalism? In *(Kolb and Mönnich, 1999)*, pages 231–258.
- Thomas Maibaum. 1974. A generalized approach to formal languages. *J. Comput. System Sci.*, 88:409–439.
- J. Mezei and Jesse Wright. 1967. Algebraic automata and contextfree sets. *Information and Control*, 11:3–29.
- Jens Michaelis, Uwe Mönnich, and Frank Morawietz. 2001. On minimalist attribute grammars and macro tree transducers. In *(Rohrer et al., 2001)*, pages 287–326.
- Uwe Mönnich. 1999. On cloning contextfreeness. In *(Kolb and Mönnich, 1999)*, pages 195–229.
- Frank Morawietz and Uwe Mönnich. 2001. A model-theoretic description of tree adjoining grammars. *ENTCS*, 53.
- Christian Rohrer, Antje Roßdeutscher, and Hans Kamp, editors. 2001. *Linguistic Form and its Computation*. University of Chicago Press.
- Stuart Shieber. 1985. Evidence against the context-freeness of natural language. *Linguistics and Philosophy*, 8:333–343.