

# Parsing and Generation as Datalog Queries

Makoto Kanazawa

National Institute of Informatics

2-1-2 Hitotsubashi, Chiyoda-ku, Tokyo, 101-8430, Japan

kanazawa@nii.ac.jp

## Abstract

We show that the problems of parsing and surface realization for grammar formalisms with “context-free” derivations, coupled with Montague semantics (under a certain restriction) can be reduced in a uniform way to Datalog query evaluation. As well as giving a polynomial-time algorithm for computing all derivation trees (in the form of a shared forest) from an input string or input logical form, this reduction has the following complexity-theoretic consequences for all such formalisms: (i) the decision problem of recognizing grammaticality (surface realizability) of an input string (logical form) is in LOGCFL; and (ii) the search problem of finding one logical form (surface string) from an input string (logical form) is in functional LOGCFL. Moreover, the generalized supplementary magic-sets rewriting of the Datalog program resulting from the reduction yields efficient Earley-style algorithms for both parsing and generation.

## 1 Introduction

The representation of context-free grammars (augmented with features) in terms of definite clause programs is well-known. In the case of a bare-bone CFG, the corresponding program is in the function-free subset of logic programming, known as *Datalog*. For example, determining whether a string John found a unicorn belongs to the language of the CFG in Figure 1 is equivalent to deciding whether the Datalog program in Figure 2 together with the *database* in (1) can derive the *query* “?- S(0, 4).”

(1) John(0, 1). found(1, 2). a(2, 3). unicorn(3, 4).

S → NP VP	V → found
VP → V NP	V → caught
V → V Conj V	Conj → and
NP → Det N	Det → a
NP → John	N → unicorn

Figure 1: A CFG.

S( <i>i</i> , <i>j</i> ) :- NP( <i>i</i> , <i>k</i> ), VP( <i>k</i> , <i>j</i> ).	V( <i>i</i> , <i>j</i> ) :- found( <i>i</i> , <i>j</i> ).
VP( <i>i</i> , <i>j</i> ) :- V( <i>i</i> , <i>k</i> ), NP( <i>k</i> , <i>j</i> ).	V( <i>i</i> , <i>j</i> ) :- caught( <i>i</i> , <i>j</i> ).
V( <i>i</i> , <i>j</i> ) :- V( <i>i</i> , <i>k</i> ), Conj( <i>k</i> , <i>l</i> ), V( <i>l</i> , <i>j</i> ).	Conj( <i>i</i> , <i>j</i> ) :- and( <i>i</i> , <i>j</i> ).
NP( <i>i</i> , <i>j</i> ) :- Det( <i>i</i> , <i>k</i> ), N( <i>k</i> , <i>j</i> ).	Det( <i>i</i> , <i>j</i> ) :- a( <i>i</i> , <i>j</i> ).
NP( <i>i</i> , <i>j</i> ) :- John( <i>i</i> , <i>j</i> ).	N( <i>i</i> , <i>j</i> ) :- unicorn( <i>i</i> , <i>j</i> ).

Figure 2: The Datalog representation of a CFG.

By *naive* (or *seminative*) bottom-up evaluation (see, e.g., Ullman, 1988), the answer to such a query can be computed in polynomial time in the size of the database for any Datalog program. By recording rule instances rather than derived facts, a packed representation of the complete set of Datalog derivation trees for a given query can also be obtained in polynomial time by the same technique. Since a Datalog derivation tree uniquely determines a grammar derivation tree, this gives a reduction of context-free recognition and parsing to query evaluation in Datalog.

In this paper, we show that a similar reduction to Datalog is possible for more powerful grammar formalisms with “context-free” derivations, such as (*multi-component*) *tree-adjoining grammars* (Joshi and Schabes, 1997; Weir, 1988), *IO macro grammars* (Fisher, 1968), and (*parallel*) *multiple context-free grammars* (Seki et al., 1991). For instance, the TAG in Figure 3 is represented by the Datalog program in Figure 4. Moreover, the method of reduc-

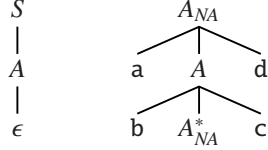


Figure 3: A TAG with one initial tree (left) and one auxiliary tree (right)

$S(p_1, p_3) :- A(p_1, p_3, p_2, p_2).$   
 $A(p_1, p_8, p_4, p_5) :- A(p_2, p_7, p_3, p_6), a(p_1, p_2), b(p_3, p_4),$   
 $c(p_5, p_6), d(p_7, p_8).$   
 $A(p_1, p_2, p_1, p_2).$

Figure 4: The Datalog representation of a TAG.

tion extends to the problem of tactical generation (surface realization) for these grammar formalisms coupled with Montague semantics (under a certain restriction). Our method essentially relies on the encoding of different formalisms in terms of *abstract categorial grammars* (de Groote, 2001).

The reduction to Datalog makes it possible to apply to parsing and generation sophisticated evaluation techniques for Datalog queries; in particular, an application of *generalized supplementary magic-sets* rewriting (Beerli and Ramakrishnan, 1991) automatically yields Earley-style algorithms for both parsing and generation. The reduction can also be used to obtain a tight upper bound, namely *LOGCFL*, on the computational complexity of the problem of recognition, both for grammaticality of input strings and for surface realizability of input logical forms.

With regard to parsing and recognition of input strings, polynomial-time algorithms and the LOGCFL upper bound on the computational complexity are already known for the grammar formalisms covered by our results (Engelfriet, 1986); nevertheless, we believe that our reduction to Datalog offers valuable insights. Concerning generation, our results seem to be entirely new.<sup>1</sup>

## 2 Context-free grammars on $\lambda$ -terms

Consider an augmentation of the grammar in Figure 1 with Montague semantics, where the left-hand

<sup>1</sup>We only consider *exact generation*, not taking into account the problem of logical form equivalence, which will most likely render the problem of generation computationally intractable (Moore, 2002).

$S(X_1 X_2) \rightarrow NP(X_1) VP(X_2)$   
 $VP(\lambda x.X_2(\lambda y.X_1 y x)) \rightarrow V(X_1) NP(X_2)$   
 $V(\lambda y x.X_2(X_1 y x)(X_3 y x)) \rightarrow V(X_1) Conj(X_2) V(X_3)$   
 $NP(X_1 X_2) \rightarrow Det(X_1) N(X_2)$   
 $NP(\lambda u.u \mathbf{John}^e) \rightarrow \mathbf{John}$   
 $V(\mathbf{find}^{e \rightarrow e \rightarrow t}) \rightarrow \mathbf{found}$   
 $V(\mathbf{catch}^{e \rightarrow e \rightarrow t}) \rightarrow \mathbf{caught}$   
 $Conj(\lambda^{t \rightarrow t \rightarrow t}) \rightarrow \mathbf{and}$   
 $Det(\lambda uv.\exists^{(e \rightarrow t) \rightarrow t}(\lambda y.\lambda^{t \rightarrow t \rightarrow t}(uy)(vy))) \rightarrow \mathbf{a}$   
 $N(\mathbf{unicorn}^{e \rightarrow t}) \rightarrow \mathbf{unicorn}$

Figure 5: A context-free grammar with Montague semantics.

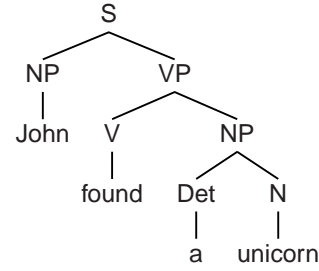


Figure 6: A derivation tree.

side of each rule is annotated with a  $\lambda$ -term that tells how the meaning of the left-hand side is composed from the meanings of the right-hand side nonterminals, represented by upper-case variables  $X_1, X_2, \dots$  (Figure 5).<sup>2</sup>

The meaning of a sentence is computed from its derivation tree. For example, *John found a unicorn* has the derivation tree in Figure 6, and the grammar rules assign its root node the  $\lambda$ -term

$(\lambda u.u \mathbf{John})(\lambda x.(\lambda uv.\exists(\lambda y.\lambda(uy)(vy)))) \mathbf{unicorn} (\lambda y.\mathbf{find} y x)$ ,

which  $\beta$ -reduces to the  $\lambda$ -term

(2)  $\exists(\lambda y.\lambda(\mathbf{unicorn} y)(\mathbf{find} y \mathbf{John}))$

encoding the first-order logic formula representing the meaning of the sentence (i.e., its logical form).

Thus, computing the logical form(s) of a sentence involves parsing and  $\lambda$ -term normalization. To find a sentence expressing a given logical form, it suffices

<sup>2</sup>We follow standard notational conventions in typed  $\lambda$ -calculus. Thus, an application  $M_1 M_2 M_3$  (written without parentheses) associates to the left,  $\lambda x.\lambda y.M$  is abbreviated to  $\lambda xy.M$ , and  $\alpha \rightarrow \beta \rightarrow \gamma$  stands for  $\alpha \rightarrow (\beta \rightarrow \gamma)$ . We refer the reader to Hindley, 1997 or Sørensen and Urzyczyn, 2006 for standard notions used in simply typed  $\lambda$ -calculus.

$S(X_1 X_2) :- NP(X_1), VP(X_2).$   
 $VP(\lambda x.X_2(\lambda y.X_1 y x)) :- V(X_1), NP(X_2).$   
 $V(\lambda y x.X_2(X_1 y x)(X_3 y x)) :- V(X_1), Conj(X_2), V(X_3).$   
 $NP(X_1 X_2) :- Det(X_1), N(X_2).$   
 $NP(\lambda u.u \text{ John}^e).$   
 $V(\text{find}^{e \rightarrow e \rightarrow t}).$   
 $V(\text{catch}^{e \rightarrow e \rightarrow t}).$   
 $Conj(\wedge^{t \rightarrow t \rightarrow t}).$   
 $Det(\lambda uv.\exists^{(e \rightarrow t) \rightarrow t}(\lambda y.\wedge^{t \rightarrow t \rightarrow t}(uy)(vy))).$   
 $N(\text{unicorn}^{e \rightarrow t}).$

Figure 7: A CFLG.

to find a derivation tree whose root node is associated with a  $\lambda$ -term that  $\beta$ -reduces to the given logical form; the desired sentence can simply be read off from the derivation tree. At the heart of both tasks is the computation of the derivation tree(s) that yield the input. In the case of generation, this may be viewed as parsing the input  $\lambda$ -term with a “context-free” grammar that generates a set of  $\lambda$ -terms (in normal form) (Figure 7), which is obtained from the original CFG with Montague semantics by stripping off terminal symbols. Determining whether a given logical form is surface realizable with the original grammar is equivalent to recognition with the resulting *context-free  $\lambda$ -term grammar* (CFLG).

In a CFLG such as in Figure 7, constants appearing in the  $\lambda$ -terms have preassigned types indicated by superscripts. There is a mapping  $\sigma$  from nonterminals to their types ( $\sigma = \{S \mapsto t, NP \mapsto (e \rightarrow t) \rightarrow t, VP \mapsto e \rightarrow t, V \mapsto e \rightarrow e \rightarrow t, Conj \mapsto t \rightarrow t \rightarrow t, Det \mapsto (e \rightarrow t) \rightarrow (e \rightarrow t) \rightarrow t, N \mapsto e \rightarrow t\}$ ). A rule that has  $A$  on the left-hand side and  $B_1, \dots, B_n$  as right-hand side nonterminals has its left-hand side annotated with a well-formed  $\lambda$ -term  $M$  that has type  $\sigma(A)$  under the type environment  $X_1 : \sigma(B_1), \dots, X_n : \sigma(B_n)$  (in symbols,  $X_1 : \sigma(B_1), \dots, X_n : \sigma(B_n) \vdash M : \sigma(A)$ ).

What we have called a context-free  $\lambda$ -term grammar is nothing but an alternative notation for an *abstract categorial grammar* (de Groote, 2001) whose abstract vocabulary is second-order, with the restriction to *linear*  $\lambda$ -terms removed.<sup>3</sup> In the linear case, Salvati (2005) has shown the recognition/parsing complexity to be PTIME, and exhibited an algorithm similar to Earley parsing for TAGs. Second-order

<sup>3</sup>A  $\lambda$ -term is a *linear*  $\lambda$ -term if each occurrence of  $\lambda$  binds at least one occurrence of a variable. A *linear*  $\lambda$ -term is *linear* if no subterm contains more than one free occurrence of the same variable.

$S(\lambda y.X_1(\lambda z.z)y) :- A(X_1).$   
 $A(\lambda xy.a^{o \rightarrow o}(X_1(\lambda z.b^{o \rightarrow o}(x(c^{o \rightarrow o}z)))(d^{o \rightarrow o}y))) :- A(X_1).$   
 $A(\lambda xy.xy).$

Figure 8: The CFLG encoding a TAG.

linear ACGs are known to be expressive enough to encode well-known mildly context-sensitive grammar formalisms in a straightforward way, including TAGs and multiple context-free grammars (de Groote, 2002; de Groote and Pogodalla, 2004).

For example, the linear CFLG in Figure 8 is an encoding of the TAG in Figure 3, where  $\sigma(S) = o \rightarrow o$  and  $\sigma(A) = (o \rightarrow o) \rightarrow o \rightarrow o$  (see de Groote, 2002 for details of this encoding). In encoding a string-generating grammar, a CFLG uses  $o$  as the type of string position and  $o \rightarrow o$  as the type of string. Each terminal symbol is represented by a constant of type  $o \rightarrow o$ , and a string  $a_1 \dots a_n$  is encoded by the  $\lambda$ -term  $\lambda z.a_1^{o \rightarrow o}(\dots(a_n^{o \rightarrow o}z)\dots)$ , which has type  $o \rightarrow o$ .

A string-generating grammar coupled with Montague semantics may be represented by a *synchronous CFLG*, a pair of CFLGs with matching rule sets (de Groote 2001). The transduction between strings and logical forms in either direction consists of parsing the input  $\lambda$ -term with the source-side grammar and normalizing the  $\lambda$ -term(s) constructed in accordance with the target-side grammar from the derivation tree(s) output by parsing.

### 3 Reduction to Datalog

We show that under a weaker condition than linearity, a CFLG can be represented by a Datalog program, obtaining a tight upper bound (LOGCFL) on the recognition complexity. Due to space limitation, our presentation here is kept at an informal level; formal definitions and rigorous proof of correctness will appear elsewhere.

We use the grammar in Figure 7 as an example, which is represented by the Datalog program in Figure 9. Note that all  $\lambda$ -terms in this grammar are *almost linear* in the sense that they are  $\lambda I$ -terms where any variable occurring free more than once in any subterm must have an atomic type. Our construction is guaranteed to be correct only when this condition is met.

Each Datalog rule is obtained from the corresponding grammar rule in the following way. Let

$S(p_1) :- NP(p_1, p_2, p_3), VP(p_2, p_3).$   
 $VP(p_1, p_4) :- V(p_2, p_4, p_3), NP(p_1, p_2, p_3).$   
 $V(p_1, p_4, p_3) :-$   
 $\quad V(p_2, p_4, p_3), Conj(p_1, p_5, p_2), V(p_5, p_4, p_3).$   
 $NP(p_1, p_4, p_5) :- Det(p_1, p_4, p_5, p_2, p_3), N(p_2, p_3).$   
 $NP(p_1, p_1, p_2) :- \mathbf{John}(p_2).$   
 $V(p_1, p_3, p_2) :- \mathbf{find}(p_1, p_3, p_2).$   
 $V(p_1, p_3, p_2) :- \mathbf{catch}(p_1, p_3, p_2).$   
 $Conj(p_1, p_3, p_2) :- \wedge(p_1, p_3, p_2).$   
 $Det(p_1, p_5, p_4, p_3, p_4) :- \exists(p_1, p_2, p_4), \wedge(p_2, p_5, p_3).$   
 $N(p_1, p_2) :- \mathbf{unicorn}(p_1, p_2).$

Figure 9: The Datalog representation of a CFLG.

$M$  be the  $\lambda$ -term annotating the left-hand side of the grammar rule. We first obtain a *principal* (i.e., most general) *typing* of  $M$ .<sup>4</sup> In the case of the second rule, this is

$$X_1 : p_3 \rightarrow p_4 \rightarrow p_2, X_2 : (p_3 \rightarrow p_2) \rightarrow p_1 \vdash \\ \lambda x.X_2(\lambda y.X_1yx) : p_4 \rightarrow p_1.$$

We then remove  $\rightarrow$  and parentheses from the types in the principal typing and write the resulting sequences of atomic types in reverse.<sup>5</sup> We obtain the Datalog rule by replacing  $X_i$  and  $M$  in the grammar rule with the sequence coming from the type paired with  $X_i$  and  $M$ , respectively. Note that atomic types in the principal typing become variables in the Datalog rule. When there are constants in the  $\lambda$ -term  $M$ , they are treated like free variables. In the case of the second-to-last rule, the principal typing is

$$\exists : (p_4 \rightarrow p_2) \rightarrow p_1, \wedge : p_3 \rightarrow p_5 \rightarrow p_2 \vdash \\ \lambda uv.\exists(\lambda y.\wedge(uy)(vy)) : (p_4 \rightarrow p_3) \rightarrow (p_4 \rightarrow p_5) \rightarrow p_1.$$

If the same constant occurs more than once, distinct occurrences are treated as distinct free variables.

The construction of the database representing the input  $\lambda$ -term is similar, but slightly more complex. A simple case is the  $\lambda$ -term (2), where each constant occurs just once. We compute its principal typing, treating constants as free variables.<sup>6</sup>

$$\exists : (4 \rightarrow 2) \rightarrow 1, \wedge : 3 \rightarrow 5 \rightarrow 2, \\ \mathbf{unicorn} : 4 \rightarrow 3, \mathbf{find} : 4 \rightarrow 6 \rightarrow 5, \mathbf{John} : 6 \\ \vdash \exists(\lambda y.\wedge(\mathbf{unicorn} y)(\mathbf{find} y \mathbf{John})) : 1.$$

<sup>4</sup>To be precise, we must first convert  $M$  to its  $\eta$ -long form relative to the type assigned to it by the grammar. For example,  $X_1 X_2$  in the first rule is converted to  $X_1(\lambda x.X_2x)$ .

<sup>5</sup>The reason for reversing the sequences of atomic types is to reconcile the  $\lambda$ -term encoding of strings with the convention of listing string positions from left to right in databases like (1).

<sup>6</sup>We assume that the input  $\lambda$ -term is in  $\eta$ -long normal form.

We then obtain the corresponding database (3) and query (4) from the antecedent and succedent of this judgment, respectively. Note that here we are using  $1, 2, 3, \dots$  as atomic types, which become database constants.

$$(3) \quad \exists(1, 2, 4). \wedge(2, 5, 3). \mathbf{unicorn}(3, 4). \\ \mathbf{find}(5, 6, 4). \mathbf{John}(6).$$

$$(4) \quad ? - S(1).$$

When the input  $\lambda$ -term contains more than one occurrence of the same constant, it is not always correct to simply treat them as distinct free variables, unlike in the case of  $\lambda$ -terms annotating grammar rules. Consider the  $\lambda$ -term (5) (John found and caught a unicorn):

$$(5) \quad \exists(\lambda y.\wedge(\mathbf{unicorn} y)(\wedge(\mathbf{find} y \mathbf{John})(\mathbf{catch} y \mathbf{John}))).$$

Here, the two occurrences of **John** must be treated as the same variable. The principal typing is (6) and the resulting database is (7).

$$(6) \quad \exists : (4 \rightarrow 2) \rightarrow 1, \wedge_1 : 3 \rightarrow 5 \rightarrow 2, \\ \mathbf{unicorn} : 4 \rightarrow 3, \wedge_2 : 6 \rightarrow 8 \rightarrow 5, \\ \mathbf{find} : 4 \rightarrow 7 \rightarrow 6, \mathbf{John} : 7, \mathbf{catch} : 4 \rightarrow 7 \rightarrow 8 \\ \vdash \exists(\lambda y.\wedge_1(\mathbf{unicorn} y) \\ (\wedge_2(\mathbf{find} y \mathbf{John})(\mathbf{catch} y \mathbf{John}))) : 1.$$

$$(7) \quad \exists(1, 2, 4). \wedge(2, 5, 3). \wedge(5, 8, 6). \mathbf{unicorn}(3, 4). \\ \mathbf{find}(6, 7, 4). \mathbf{John}(7). \mathbf{catch}(8, 7, 4).$$

It is not correct to identify the two occurrences of  $\wedge$  in this example. The rule is to identify distinct occurrences of the same constant just in case they occur in the same position within  $\alpha$ -equivalent subterms of an atomic type. This is a necessary condition for those occurrences to originate as one and the same occurrence in the non-normal  $\lambda$ -term at the root of the derivation tree. (As a preprocessing step, it is also necessary to check that distinct occurrences of a bound variable satisfy the same condition, so that the given  $\lambda$ -term is  $\beta$ -equal to some almost linear  $\lambda$ -term.<sup>7</sup>)

<sup>7</sup>Note that the way we obtain a database from an input  $\lambda$ -term generalizes the standard database representation of a string: from the  $\lambda$ -term encoding  $\lambda z.a_1^{o \rightarrow o}(\dots(a_n^{o \rightarrow o} z)\dots)$  of a string  $a_1 \dots a_n$ , we obtain the database  $\{a_1(0, 1), \dots, a_n(n-1, n)\}$ .

## 4 Correctness of the reduction

We sketch some key points in the proof of correctness of our reduction. The  $\lambda$ -term  $N$  obtained from the input  $\lambda$ -term by replacing occurrences of constants by free variables in the manner described above is the normal form of some almost linear  $\lambda$ -term  $N'$ . The *leftmost reduction* from an almost linear  $\lambda$ -term to its normal form must be *non-deleting* and *almost non-duplicating* in the sense that when a  $\beta$ -redex  $(\lambda x.P)Q$  is contracted,  $Q$  is not deleted, and moreover it is not duplicated unless the type of  $x$  is atomic. We can show that the *Subject Expansion Theorem* holds for such  $\beta$ -reduction, so the principal typing of  $N$  is also the principal typing of  $N'$ . By a slight generalization of a result by Aoto (1999), this typing  $\Gamma \vdash N' : \alpha$  must be *negatively non-duplicated* in the sense that each atomic type has at most one negative occurrence in it. By Aoto and Ono's (1994) generalization of the *Coherence Theorem* (see Mints, 2000), it follows that every  $\lambda$ -term  $P$  such that  $\Gamma' \vdash P : \alpha$  for some  $\Gamma' \subseteq \Gamma$  must be  $\beta\eta$ -equal to  $N'$  (and consequently to  $N$ ).

Given the one-one correspondence between the grammar rules and the Datalog rules, a Datalog derivation tree uniquely determines a grammar derivation tree (see Figure 10 as an example). This relation is not one-one, because a Datalog derivation tree contains database constants from the input database. This extra information determines a typing of the  $\lambda$ -term  $P$  at the root of the grammar derivation tree (with occurrences of constants in the  $\lambda$ -term corresponding to distinct facts in the database regarded as distinct free variables):

**John** : 6, **find** :  $4 \rightarrow 6 \rightarrow 5$ ,  $\exists : (4 \rightarrow 2) \rightarrow 1$ ,  
 $\wedge : 3 \rightarrow 5 \rightarrow 2$ , **unicorn** :  $4 \rightarrow 3 \vdash$   
 $(\lambda u.u \text{ John})$   
 $(\lambda x.(\lambda uv.\exists(\lambda y.\wedge(uv)(vy)))) \text{ unicorn } (\lambda y.\text{find } y \ x) : 1.$

The antecedent of this typing must be a subset of the antecedent of the principal typing of the  $\lambda$ -term  $N$  from which the input database was obtained. By the property mentioned at the end of the preceding paragraph, it follows that the grammar derivation tree is a derivation tree for the input  $\lambda$ -term.

Conversely, consider the  $\lambda$ -term  $P$  (with distinct occurrences of constants regarded as distinct free variables) at the root of a grammar derivation tree

for the input  $\lambda$ -term. We can show that there is a substitution  $\theta$  which maps the free variables of  $P$  to the free variables of the  $\lambda$ -term  $N$  used to build the input database such that  $\theta$  sends the normal form of  $P$  to  $N$ . Since  $P$  is an almost linear  $\lambda$ -term, the leftmost reduction from  $P\theta$  to  $N$  is non-deleting and almost non-duplicating. By the Subject Expansion Theorem, the principal typing of  $N$  is also the principal typing of  $P\theta$ , and this together with the grammar derivation tree determines a Datalog derivation tree.

## 5 Complexity-theoretic consequences

Let us call a rule  $A(M) :- B_1(X_1), \dots, B_n(X_n)$  in a CFLG an  $\epsilon$ -rule if  $n = 0$  and  $M$  does not contain any constants. We can eliminate  $\epsilon$ -rules from an almost linear CFLG by the same method that Kanazawa and Yoshinaka (2005) used for linear grammars, noting that for any  $\Gamma$  and  $\alpha$ , there are only finitely many almost linear  $\lambda$ -terms  $M$  such that  $\Gamma \vdash M : \alpha$ . If a grammar has no  $\epsilon$ -rule, any derivation tree for the input  $\lambda$ -term  $N$  that has a  $\lambda$ -term  $P$  at its root node corresponds to a Datalog derivation tree whose number of leaves is equal to the number of occurrences of constants in  $P$ , which cannot exceed the number of occurrences of constants in  $N$ .

A Datalog program  $\mathbf{P}$  is said to have the *polynomial fringe property* relative to a class  $\mathcal{D}$  of databases if there is a polynomial  $p(n)$  such that for every database  $D$  in  $\mathcal{D}$  of  $n$  facts and every query  $q$  such that  $\mathbf{P} \cup D$  derives  $q$ , there is a derivation tree for  $q$  whose fringe (i.e., sequence of leaves) is of length at most  $p(n)$ . For such  $\mathbf{P}$  and  $\mathcal{D}$ , it is known that  $\{(D, q) \mid D \in \mathcal{D}, \mathbf{P} \cup D \text{ derives } q\}$  is in the complexity class *LOGCFL* (Ullman and Van Gelder, 1988; Kanellakis, 1988).

We state without proof that the database-query pair  $(D, q)$  representing an input  $\lambda$ -term  $N$  can be computed in logspace. By padding  $D$  with extra useless facts so that the size of  $D$  becomes equal to the number of occurrences of constants in  $N$ , we obtain a logspace reduction from the set of  $\lambda$ -terms generated by an almost linear CFLG to a set of the form  $\{(D, q) \mid D \in \mathcal{D}, \mathbf{P} \cup D \vdash q\}$ , where  $\mathbf{P}$  has the polynomial fringe property relative to  $\mathcal{D}$ . This shows that the problem of recognition for an almost linear CFLG is in LOGCFL.



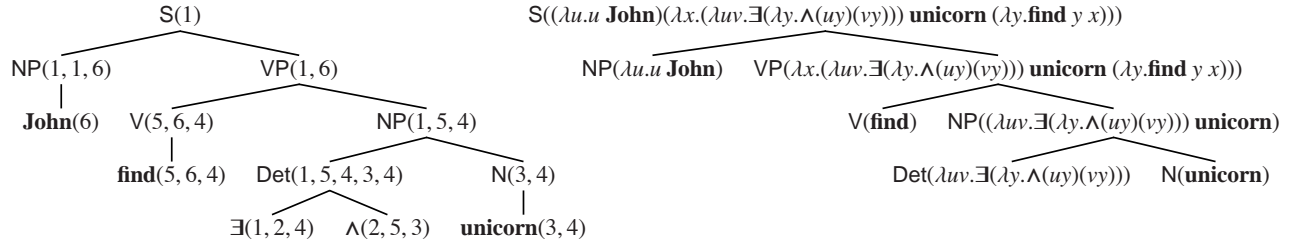


Figure 10: A Datalog derivation tree (left) and the corresponding grammar derivation tree (right)

By the main result of Gottlob et al. (2002), the related search problem of finding one derivation tree for the input  $\lambda$ -term is in *functional LOGCFL*, i.e., the class of functions that can be computed by a logspace-bounded Turing machine with a LOGCFL oracle. In the case of a synchronous almost linear CFLG, the derivation tree found from the source  $\lambda$ -term can be used to compute a target  $\lambda$ -term. Thus, to the extent that transduction back and forth between strings and logical forms can be expressed by a synchronous almost linear CFLG, the search problem of finding one logical form of an input sentence and that of finding one surface realization of an input logical form are both in functional LOGCFL.<sup>8</sup> As a consequence, there are efficient parallel algorithms for these problems.

## 6 Regular sets of trees as input

Almost linear CFLGs can represent a substantial fragment of a Montague semantics for English and such “linear” grammar formalisms as (multi-component) tree-adjoining grammars (both as string grammars and as tree grammars) and multiple context-free grammars. However, IO macro grammars and parallel multiple context-free grammars cannot be directly represented because representing string copying requires multiple occurrences of a variable of type  $o \rightarrow o$ . This problem can be solved by switching from strings to trees. We convert the input string into the regular set of binary trees whose yield equals the input string (using  $c$

<sup>8</sup>If the target-side grammar is not linear, the normal form of the target  $\lambda$ -term cannot be explicitly computed because its size may be exponential in the size of the source  $\lambda$ -term. Nevertheless, a typing that serves to uniquely identify the target  $\lambda$ -term can be computed from the derivation tree in logspace. Also, if the target-side grammar is linear and string-generating, the target string can be explicitly computed from the derivation tree in logspace (Salvati, 2007).

as the sole symbol of rank 2), and turn the grammar into a tree grammar, replacing all instances of string concatenation in the grammar with the tree operation  $t_1, t_2 \mapsto c(t_1, t_2)$ . This way, a string grammar is turned into a tree grammar that generates a set of trees whose image under the yield function is the language of the string grammar. (In the case of an IO macro grammar, the result is an *IO context-free tree grammar* (Engelfriet, 1977).) String copying becomes tree copying, and the resulting grammar can be represented by an almost linear CFLG and hence by a Datalog program. The regular set of all binary trees that yield the input string is represented by a database that is constructed from a deterministic bottom-up finite tree automaton recognizing it. Determinism is important for ensuring correctness of this reduction. Since the database can be computed from the input string in logspace, the complexity-theoretic consequences of the last section carry over here.

## 7 Magic sets and Earley-style algorithms

*Magic-sets* rewriting of a Datalog program allows bottom-up evaluation to avoid deriving useless facts by mimicking top-down evaluation of the original program. The result of the *generalized supplementary magic-sets* rewriting of Beerli and Ramakrishnan (1991) applied to the Datalog program representing a CFG essentially coincides with the *deduction system* (Shieber et al., 1995) or *uninstantiated parsing system* (Sikkel, 1997) for Earley parsing. By applying the same rewriting method to Datalog programs representing almost linear CFLGs, we can obtain efficient parsing and generation algorithms for various grammar formalisms with context-free derivations.

We illustrate this approach with the program in Figure 4, following the presentation of Ullman

(1989a; 1989b). We assume the query to take the form “ $?- S(0, x)$ ”, so that the input database can be processed incrementally. The program is first made *safe* by eliminating the possibility of deriving non-ground atoms:

$$\begin{aligned} S(p_1, p_3) &:- A(p_1, p_3, p_2, p_2). \\ A(p_1, p_8, p_4, p_5) &:- A(p_2, p_7, p_3, p_6), a(p_1, p_2), b(p_3, p_4), c(p_5, p_6), d(p_7, p_8). \\ A(p_1, p_8, p_4, p_5) &:- a(p_1, p_2), b(p_2, p_4), c(p_5, p_6), d(p_6, p_8). \end{aligned}$$

The *subgoal rectification* removes duplicate arguments from subgoals, creating new predicates as needed:

$$\begin{aligned} S(p_1, p_3) &:- B(p_1, p_3, p_2). \\ A(p_1, p_8, p_4, p_5) &:- A(p_2, p_7, p_3, p_6), a(p_1, p_2), b(p_3, p_4), c(p_5, p_6), d(p_7, p_8). \\ A(p_1, p_8, p_4, p_5) &:- a(p_1, p_2), b(p_2, p_4), c(p_5, p_6), d(p_6, p_8). \\ B(p_1, p_8, p_4) &:- A(p_2, p_7, p_3, p_6), a(p_1, p_2), b(p_3, p_4), c(p_4, p_6), d(p_7, p_8). \\ B(p_1, p_8, p_4) &:- a(p_1, p_2), b(p_2, p_4), c(p_4, p_6), d(p_6, p_8). \end{aligned}$$

We then attach to predicates *adornments* indicating the free/bound status of arguments in top-down evaluation, reordering subgoals so that as many arguments as possible are marked as bound:

$$\begin{aligned} S^{bf}(p_1, p_3) &:- B^{bff}(p_1, p_3, p_2). \\ B^{bff}(p_1, p_8, p_4) &:- a^{bf}(p_1, p_2), A^{bff}(p_2, p_7, p_3, p_6), b^{bf}(p_3, p_4), c^{bb}(p_4, p_6), \\ &\quad d^{bf}(p_7, p_8). \\ B^{bff}(p_1, p_8, p_4) &:- a^{bf}(p_1, p_2), b^{bf}(p_2, p_4), c^{bf}(p_4, p_6), d^{bf}(p_6, p_8). \\ A^{bff}(p_1, p_8, p_4, p_5) &:- a^{bf}(p_1, p_2), A^{bff}(p_2, p_7, p_3, p_6), b^{bf}(p_3, p_4), c^{bb}(p_5, p_6), \\ &\quad d^{bf}(p_7, p_8). \\ A^{bff}(p_1, p_8, p_4, p_5) &:- a^{bf}(p_1, p_2), b^{bf}(p_2, p_4), c^{ff}(p_5, p_6), d^{bf}(p_6, p_8). \end{aligned}$$

The generalized supplementary magic-sets rewriting finally gives the following rule set:

$$\begin{aligned} r_1 &: m\_B(p_1) :- m\_S(p_1). \\ r_2 &: S(p_1, p_3) :- m\_B(p_1), B(p_1, p_3, p_2). \\ r_3 &: sup_{2,1}(p_1, p_2) :- m\_B(p_1), a(p_1, p_2). \\ r_4 &: sup_{2,2}(p_1, p_7, p_3, p_6) :- sup_{2,1}(p_1, p_2), A(p_2, p_7, p_3, p_6). \\ r_5 &: sup_{2,3}(p_1, p_7, p_6, p_4) :- sup_{2,2}(p_1, p_7, p_3, p_6), b(p_3, p_4). \\ r_6 &: sup_{2,4}(p_1, p_7, p_4) :- sup_{2,3}(p_1, p_7, p_6, p_4), c(p_4, p_6). \\ r_7 &: B(p_1, p_8, p_4) :- sup_{2,4}(p_1, p_7, p_4), d(p_7, p_8). \\ r_8 &: sup_{3,1}(p_1, p_2) :- m\_B(p_1), a(p_1, p_2). \\ r_9 &: sup_{3,2}(p_1, p_4) :- sup_{3,1}(p_1, p_2), b(p_2, p_4). \\ r_{10} &: sup_{3,3}(p_1, p_4, p_6) :- sup_{3,2}(p_1, p_4), c(p_4, p_6). \\ r_{11} &: B(p_1, p_8, p_4) :- sup_{3,3}(p_1, p_4, p_6), d(p_6, p_8). \\ r_{12} &: m\_A(p_2) :- sup_{2,1}(p_1, p_2). \\ r_{13} &: m\_A(p_2) :- sup_{4,1}(p_1, p_2). \\ r_{14} &: sup_{4,1}(p_1, p_2) :- m\_A(p_1), a(p_1, p_2). \\ r_{15} &: sup_{4,2}(p_1, p_7, p_3, p_6) :- sup_{4,1}(p_1, p_2), A(p_2, p_7, p_3, p_6). \\ r_{16} &: sup_{4,3}(p_1, p_7, p_6, p_4) :- sup_{4,2}(p_1, p_7, p_3, p_6), b(p_3, p_4). \\ r_{17} &: sup_{4,4}(p_1, p_7, p_4, p_5) :- sup_{4,3}(p_1, p_7, p_6, p_4), c(p_5, p_6). \\ r_{18} &: A(p_1, p_8, p_4, p_5) :- sup_{4,4}(p_1, p_7, p_4, p_5), d(p_7, p_8). \\ r_{19} &: sup_{5,1}(p_1, p_2) :- m\_A(p_1), a(p_1, p_2). \\ r_{20} &: sup_{5,2}(p_1, p_4) :- sup_{5,1}(p_1, p_2), b(p_2, p_4). \\ r_{21} &: sup_{5,3}(p_1, p_4, p_5, p_6) :- sup_{5,2}(p_1, p_4), c(p_5, p_6). \\ r_{22} &: A(p_1, p_8, p_4, p_5) :- sup_{5,3}(p_1, p_4, p_5, p_6), d(p_6, p_8). \end{aligned}$$

The following version of chart parsing adds control structure to this deduction system:

1. (INIT) Initialize the chart to the empty set, the agenda to the singleton  $\{m\_S(0)\}$ , and  $n$  to 0.
2. Repeat the following steps:
  - (a) Repeat the following steps until the agenda is exhausted:
    - i. Remove a fact from the agenda, called the *trigger*.
    - ii. Add the trigger to the chart.
    - iii. Generate all facts that are immediate consequences of the trigger together with all facts in the chart, and add to the agenda those generated facts that are neither already in the chart nor in the agenda.
  - (b) (SCAN) Remove the next fact from the input database and add it to the agenda, incrementing  $n$ . If there is no more fact in the input database, go to step 3.
3. If  $S(0, n)$  is in the chart, accept; otherwise reject.

The following is the trace of the algorithm on input string aabbccdd:

1. $m\_S(0)$	INIT	14. $c(4, 5)$	SCAN
2. $m\_B(0)$	$r_1, 1$	15. $sup_{5,3}(1, 3, 4, 5)$	$r_{21}, 12, 14$
3. $a(0, 1)$	SCAN	16. $c(6, 5)$	SCAN
4. $sup_{2,1}(0, 1)$	$r_3, 2, 3$	17. $sup_{5,3}(1, 3, 5, 6)$	$r_{21}, 12, 16$
5. $sup_{3,1}(0, 1)$	$r_8, 2, 3$	18. $d(6, 7)$	SCAN
6. $m\_A(1)$	$r_{12}, 4$	19. $A(1, 7, 3, 5)$	$r_{22}, 17, 18$
7. $a(1, 2)$	SCAN	20. $sup_{2,2}(0, 7, 3, 5)$	$r_4, 4, 19$
8. $sup_{4,1}(1, 2)$	$r_{14}, 6, 7$	21. $sup_{2,3}(0, 7, 5, 4)$	$r_5, 13, 20$
9. $sup_{5,1}(1, 2)$	$r_{19}, 6, 7$	22. $sup_{2,4}(0, 7, 4)$	$r_6, 14, 21$
10. $m\_A(2)$	$r_{13}, 8$	23. $d(7, 8)$	SCAN
11. $b(2, 3)$	SCAN	24. $B(0, 8, 4)$	$r_7, 22, 23$
12. $sup_{5,2}(1, 3)$	$r_{20}, 9, 11$	25. $S(0, 8)$	$r_2, 2, 24$
13. $b(3, 4)$	SCAN		

Note that unlike existing Earley-style parsing algorithms for TAGs, the present algorithm is an instantiation of a general schema that applies to parsing with more powerful grammar formalisms as well as to generation with Montague semantics.

## 8 Conclusion

Our reduction to Datalog brings sophisticated techniques for Datalog query evaluation to the problems

of parsing and generation, and establishes a tight bound on the computational complexity of recognition for a wide range of grammars. In particular, it shows that the use of higher-order  $\lambda$ -terms for semantic representation need not be avoided for the purpose of achieving computational tractability.

## References

- Aoto, Takahito. 1999. Uniqueness of normal proofs in implicational intuitionistic logic. *Journal of Logic, Language and Information* **8**, 217–242.
- Aoto, Takahito and Hiroakira Ono. 1994. Uniqueness of normal proofs in  $\{\rightarrow, \wedge\}$ -fragment of NJ. Research Report IS-RR-94-0024F. School of Information Science, Japan Advanced Institute of Science and Technology.
- Beeri, Catriel and Raghu Ramakrishnan. 1991. On the power of magic. *Journal of Logic Programming* **10**, 255–299.
- Engelfriet, J. and E. M. Schmidt. 1977. IO and OI, part I. *The Journal of Computer and System Sciences* **15**, 328–353.
- Engelfriet, Joost. 1986. The complexity of languages generated by attribute grammars. *SIAM Journal on Computing* **15**, 70–86.
- Fisher, Michael J. 1968. *Grammars with Macro-Like Productions*. Ph.D. dissertation. Harvard University.
- Gottlob, Georg, Nicola Lenoe, Francesco Scarcello. 2002. Computing LOGCFL certificates. *Theoretical Computer Science* **270**, 761–777.
- de Groote, Philippe. 2001. Towards abstract categorial grammars. In *Association for Computational Linguistics, 39th Annual Meeting and 10th Conference of the European Chapter, Proceedings of the Conference*, pages 148–155.
- de Groote, Philippe. 2002. Tree-adjointing grammars as abstract categorial grammars. In *Proceedings of the Sixth International Workshop on Tree Adjoining Grammar and Related Frameworks (TAG+6)*, pages 145–150. Università di Venezia.
- de Groote, Philippe and Sylvain Pogodalla. 2004. On the expressive power of abstract categorial grammars: Representing context-free formalisms. *Journal of Logic, Language and Information* **13**, 421–438.
- Hindley, J. Roger. 1997. *Basic Simple Type Theory*. Cambridge: Cambridge University Press.
- Aravind K. Joshi and Yves Schabes. 1997. Tree-adjointing grammars. In Grzegorz Rozenberg and Arto Salomaa, editors, *Handbook of Formal Languages, Vol. 3*, pages 69–123. Berlin: Springer.
- Kanazawa, Makoto and Ryo Yoshinaka. 2005. Lexicalization of second-order ACGs. NII Technical Report. NII-2005-012E. National Institute of Informatics, Tokyo.
- Kanellakis, Paris C. 1988. Logic programming and parallel complexity. In Jack Minker, editor, *Foundations of Deductive Databases and Logic Programming*, pages 547–585. Los Altos, CA: Morgan Kaufmann.
- Mints, Grigori. 2000. *A Short Introduction to Intuitionistic Logic*. New York: Kluwer Academic/Plenum Publishers.
- Moore, Robert C. 2002. A complete, efficient sentence-realization algorithm for unification grammar. In *Proceedings, International Natural Language Generation Conference*, Harriman, New York, pages 41–48.
- Salvati, Sylvain. 2005. *Problèmes de filtrage et problèmes d’analyse pour les grammaires catégorielles abstraites*. Doctoral dissertation, l’Institut National Polytechnique de Lorraine.
- Salvati, Sylvain. 2007. Encoding second order string ACG with deterministic tree walking transducers. In Shuly Wintner, editor, *Proceedings of FG 2006: The 11th conference on Formal Grammar*, pages 143–156. FG Online Proceedings. Stanford, CA: CSLI Publications.
- Seki, Hiroyuki, Takashi Matsumura, Mamoru Fujii, and Tadao Kasami. 1991. On multiple context-free grammars. *Theoretical Computer Science* **88**, 191–229.
- Shieber, Stuart M., Yves Schabes, and Fernando C. N. Pereira. 1995. Principles and implementations of deductive parsing. *Journal of Logic Programming* **24**, 3–36.
- Sikkel, Klaas. 1997. *Parsing Schemata*. Berlin: Springer.
- Sørensen, Morten Heine and Paweł Urzyczyn. 2006. *Lectures on the Curry-Howard Isomorphism*. Amsterdam: Elsevier.
- Ullman, Jeffrey D. 1988. *Principles of Database and Knowledge-Base Systems. Volume I*. Rockville, MD.: Computer Science Press.
- Ullman, Jeffrey D. 1989a. Bottom-up beats top-down for Datalog. In *Proceedings of the Eighth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems*, Philadelphia, pages 140–149.
- Ullman, Jeffrey D. 1989b. *Principles of Database and Knowledge-Base Systems. Volume II: The New Technologies*. Rockville, MD.: Computer Science Press.
- Ullman, Jeffrey D. and Allen Van Gelder. 1988. Parallel complexity of logical query programs. *Algorithmica* **3**, 5–42.
- David J. Weir. 1988. *Characterizing Mildly Context-Sensitive Grammar Formalisms*. Ph.D. dissertation. University of Pennsylvania.