

# On Computing all Abductive Explanations

Thomas Eiter

Institut für Informationssysteme,  
Technische Universität Wien  
Favoritenstraße 9–11, A-1040 Wien, Austria  
e-mail: eiter@kr.tuwien.ac.at

Kazuhisa Makino

Division of Systems Science,  
Graduate School of Engineering Science,  
Osaka University, Toyonaka, Osaka 560-8531  
makino@sys.es.osaka-u.ac.jp

## Abstract

We consider the computation of all respectively a polynomial subset of the explanations of an abductive query from a Horn theory, and pay particular attention to whether the query is a positive or negative letter, the explanation is based on literals from an assumption set, and the Horn theory is represented in terms of formulas or characteristic models. We derive tractability results, one of which refutes a conjecture by Selman and Levesque, as well as intractability results, and furthermore also semi-tractability results in terms of solvability in quasi-polynomial time. Our results complement previous results in the literature, and elucidate the computational complexity of generating the set of explanations.

## Introduction

Abduction is a fundamental mode of reasoning, which has been recognized as an important principle of common-sense reasoning (see e.g. (Brewka, Dix, & Konolige 1997)). It has applications in many areas of AI including diagnosis, planning, learning, natural language understanding and many others (see e.g. references in (Eiter & Gottlob 1995)). In a logic-based setting, abduction can be defined as the task, given a set of formulas  $\Sigma$  (the background theory) and a formula  $\chi$  (the query), to find a smallest set of formulas  $E$  (an explanation) from a set of hypotheses such that  $\Sigma$  plus  $E$  is satisfiable and logically entails  $\chi$ . For use in practice, the computation of abductive explanations is an important problem, for which well-known early systems such as Theorist (Poole 1989) or ATMS solvers have been devised. Since then, there has been a vastly growing literature on this subject, indicating the need for efficient abductive procedures.

**Main problems considered.** In this paper, we consider computing a set of explanations for queries from Horn theories. More precisely, we address the following problems:

- Computing *all* explanations of a query  $\chi$  given by a letter  $q$ , with and without a set of assumption literals  $A$  from which explanations  $E$  must be formed, similar as in (Poole 1989; Selman & Levesque 1996). Note that the logical disjunction of all explanations is a weakest disjunctive form over the hypotheses explaining  $\chi$ . It is easy to see that in general, there might be exponentially many explanations, and computing all explanations is inevitably exponential. However,

Copyright © 2002, American Association for Artificial Intelligence (www.aaai.org). All rights reserved.

it is in this case of interest whether the computation is possible in *polynomial total time* (or *output-polynomial time*), i.e., in time polynomial in the combined size of the input and the output. Furthermore, if exponential space is prohibitive, it is of interest to know whether a few explanations (e.g., polynomially many) can be generated in polynomial time, as studied by Selman and Levesque (1996).

- We contrast formula-based (syntactic) with model-based (semantic) representation of Horn theories. The latter form of representation, where a Horn theory is represented by the *characteristic models*, was advocated by Kautz *et al.* (1993). As they showed, important inference problems are tractable in the model-based setting. Namely, whether a Horn theory  $\Sigma$  implies a CNF  $\varphi$ , and whether a query  $q$  has an explanation w.r.t. an assumption set  $A$ ; the latter is intractable under formula-based representation. Similar results were shown for other theories by Khardon and Roth (1996).

- We investigate the role of syntax for computing abductive queries. In the framework of (Selman & Levesque 1996; Kautz, Kearns, & Selman 1993), the query is a positive letter  $q$ . However, it is of equal interest to consider *negative queries* as well, i.e., to explain the complement  $\bar{q}$  of an atom  $q$ . Since the Horn property imposes semantic restrictions on theories, it is not straightforward to express such negative queries in terms of positive queries.

- Finally, we consider as a meaningful generalization the computation of *joint explanations*. That is, given a background theory  $\Sigma$  and observations  $o_1, o_2, \dots, o_l$ , where  $l \geq 2$ , compute a *single* explanation  $E$  which is good for *each*  $o_i$ . Joint explanations are relevant, e.g., in diagnostic reasoning. We may want to know whether different observations allow to come up with the same diagnosis, given by an abductive explanation, about a system malfunctioning. Such a diagnosis is particularly strong, as it is backed up by several cases.

**Main results.** Our main results are summarized as follows.

- We refute Selman and Levesque's belief (1996, p. 266), that given a Horn theory  $\Sigma$  and a query letter  $q$ , it is hard to list all explanations of  $q$  from  $\Sigma$  even if we are *guaranteed* that there are only few explanations. More precisely, we disprove their conjecture that generating  $O(n)$  many explanations of  $q$  is NP-hard, where  $n$  is the number of propositional letters in the language. This is a consequence of our result that

generating *all* nontrivial explanations of  $q$  is possible in *total polynomial time* (Theorem 1).

- We give a detailed characterization of computing all explanations of a query from a Horn theory in the formula- vs model-based setting, for both general explanations and explanations w.r.t. a set of assumption literals. In a nutshell, we obtain three kinds of results:

(1) A procedure which enumerates all nontrivial explanations of a query letter  $q$  from a Horn theory  $\Sigma$  with incremental polynomial delay. This is a positive result and trivially implies that all explanations can be found in polynomial total time. Moreover, it means that any polynomial number of explanations can be generated in polynomial time in the size of the input (Corollary 1).

(2) Intractability results for generating all explanations for a negative query  $\bar{q}$  from a Horn theory  $\Sigma$  contained in a set of assumption literals  $A$ ; this complements a similar result for positive queries in (Selman & Levesque 1996). Both results emerge from the fact that the associated problems of recognizing the correct output are co-NP-complete. Since some hard instances have only small (polynomial-size) output, they also imply that computing few (polynomially many) explanations is intractable.

(3) Under model-based representation, generating all explanations is polynomial-time equivalent to the well-known problem of dualizing a positive CNF (DUALIZATION), i.e., given a CNF  $\varphi$  in which no negative literal occurs, compute the (unique) prime DNF of  $\varphi$ . DUALIZATION is a well-known open problem in NP-completeness, cf. (Bioch & Ibaraki 1995; Fredman & Khachiyan 1996); it is known to be solvable in *quasi-polynomial* total time, i.e., in time  $N^{o(\log N)}$  where  $N$  denotes the combined size of the input and output (Fredman & Khachiyan 1996); furthermore, polynomial total time algorithms are known for many special cases, and as recently shown, the decisional variant of recognizing the prime DNF of  $\varphi$  is solvable with limited nondeterminism, i.e., in polynomial time with  $O(\log^2 N)$  many guesses (Eiter, Gottlob, & Makino 2002). Since DUALIZATION is strongly believed not to be co-NP-hard, our result thus provides strong evidence that under model-based representation, computing all explanations is not co-NP-hard. Interestingly, the equivalence result holds for both positive and negative queries, and whether arbitrary or explanations over a set of assumption literals  $A$  are admitted. This means that, in a sense, model-based representation, in contrast to formula-based representation, is not sensitive to these aspects. Furthermore, by resorting to respective algorithms for dualization, the result provides us with algorithms for enumerating all or polynomially many explanations with quasi-polynomial time delay between outputs.

- We show that deciding the existence of a joint explanation is intractable, for both formula- and model-based representation. Thus, the positive results for ordinary explanations do not extend to joint explanations.

Proofs of all results are given in (Eiter & Makino 2002).

## Preliminaries

We assume a standard propositional language with letters  $x_1, x_2, \dots, x_n$  from a set  $P$ , where each  $x_i$  takes either value

1 (true) or 0 (false). Negated atoms are denoted by  $\bar{x}_i$ , and the opposite of a literal  $\ell$  by  $\bar{\ell}$ . Furthermore, we use  $\bar{A} = \{\bar{\ell} \mid \ell \in A\}$  for any set of literals  $A$  and set  $Lit = P \cup \bar{P}$ .

A clause is a disjunction  $c = \bigvee_{p \in P(c)} p \vee \bigvee_{p \in N(c)} \bar{p}$  of literals, where  $P(c)$  and  $N(c)$  are the sets of atoms occurring positive and negated in  $c$  and  $P(c) \cap N(c) = \emptyset$ . Dually, a term is conjunction  $t = \bigwedge_{p \in P(t)} p \wedge \bigwedge_{p \in N(t)} \bar{p}$  of literals, where  $P(t)$  and  $N(t)$  are similarly defined. We also view clauses and terms as sets of literals  $P(c) \cup N(c)$  and  $P(t) \cup N(t)$ , respectively. A clause  $c$  is *Horn*, if  $|P(c)| \leq 1$ , and a CNF is *Horn*, if it contains only Horn clauses. A theory  $\Sigma$  is any set of formulas; it is *Horn*, if it is a set of Horn clauses. As usual, we identify  $\Sigma$  with  $\varphi = \bigwedge_{c \in \Sigma} c$ , and write  $c \in \varphi$  etc.

**Definition 1** Given a (Horn) theory  $\Sigma$ , called the background theory, a letter  $q$  (called query), and a set of literals  $A \subseteq Lit$ , an *explanation of  $q$  w.r.t.  $A$*  is a minimal set of literals  $E$  over  $A$  such that

- (i)  $\Sigma \cup E \models q$ , and
- (ii)  $\Sigma \cup E$  is satisfiable.

If  $A = Lit$ , then we call  $E$  simply an *explanation of  $q$* .

Observe that the above definition is slightly more general than the *assumption-based explanations* of (Selman & Levesque 1996), which emerge as  $A = P' \cup \bar{P}'$  where  $P' \subseteq P$  (i.e.,  $A$  contains all literals over a subset  $P'$  of the letters). Furthermore, in some texts explanations must be sets of positive literals. As for Horn theories, the following is known, cf. (Khardon & Roth 1996):

**Proposition 1** Let  $E$  be any explanation of  $q$  w.r.t.  $A \subseteq Lit$ . Then  $E \subseteq P$ , i.e.,  $E$  contains only positive literals.

**Example 1** Consider a theory  $\Sigma = \{\bar{x}_1 \vee \bar{x}_4, \bar{x}_4 \vee \bar{x}_3, \bar{x}_1 \vee x_2, \bar{x}_4 \vee \bar{x}_5 \vee x_1\}$ . Suppose we want to explain  $q = x_2$  from  $A = \{x_1, x_4\}$ . Then, we find that  $E = \{x_1\}$  is an explanation. Indeed,  $\Sigma \cup \{x_1\} \models x_2$ , and  $\Sigma \cup \{x_1\}$  is satisfiable; moreover,  $E$  is minimal. On the other hand,  $E' = \{x_1, \bar{x}_4\}$  satisfies (i) and (ii) for  $q = x_2$ , but is not minimal.  $\square$

Horn theories have a well-known semantic characterization. A *model* is a vector  $v \in \{0, 1\}^n$ , whose  $i$ -th component is denoted by  $v_i$ . For  $B \subseteq \{1, \dots, n\}$ , we let  $x^B$  be the model  $v$  such that  $v_i = 1$ , if  $i \in B$  and  $v_i = 0$ , if  $i \notin B$ , for  $i \in \{1, \dots, n\}$ . The set of models of formula  $\varphi$  (resp. theory  $\Sigma$ ), denoted by  $mod(\varphi)$  (resp.  $mod(\Sigma)$ ), is defined as usual.

For models  $v, w$ , we denote by  $v \leq w$  the usual componentwise ordering, i.e.,  $v_i \leq w_i$  for all  $i = 1, 2, \dots, n$ , where  $0 \leq 1$ ;  $v < w$  means  $v \neq w$  and  $v \leq w$ . For any set of models  $\mathcal{M}$ , we denote by  $\max(\mathcal{M})$ , (resp.,  $\min(\mathcal{M})$ ) the set of all maximal (resp., minimal) models in  $\mathcal{M}$ . Denote by  $v \wedge w$  componentwise AND of vectors  $v, w \in \{0, 1\}^n$ , and by  $Cl_\wedge(S)$  the closure of  $S \subseteq \{0, 1\}^n$  under  $\wedge$ . Then, a theory  $\Sigma$  is Horn representable, iff  $mod(\Sigma) = Cl_\wedge(mod(\Sigma))$ .

**Example 2** Consider  $\mathcal{M}_1 = \{(0101), (1001), (1000)\}$  and  $\mathcal{M}_2 = \{(0101), (1001), (1000), (0001), (0000)\}$ . Then, for  $v = (0101)$ ,  $w = (1000)$ , we have  $w, v \in \mathcal{M}_1$ , while  $v \wedge w = (0000) \notin \mathcal{M}_1$ ; hence  $\mathcal{M}_1$  is not the set of models of a Horn theory. On the other hand,  $Cl_\wedge(\mathcal{M}_2) = \mathcal{M}_2$ , thus  $\mathcal{M}_2 = mod(\Sigma)$  for some Horn theory  $\Sigma$ .

As discussed by Kautz *et al.* (1993), a Horn theory  $\Sigma$  is semantically represented by its characteristic models, where  $v \in \text{mod}(\Sigma)$  is called *characteristic* (or *extreme* (Dechter & Pearl 1992)), if  $v \notin \text{Cl}_\wedge(\text{mod}(\Sigma) \setminus \{v\})$ . The set of all such models, the *characteristic set* of  $\Sigma$ , is denoted by  $\text{char}(\Sigma)$ . Note that  $\text{char}(\Sigma)$  is unique. E.g.,  $(0101) \in \text{char}(\Sigma_2)$ , while  $(0000) \notin \text{char}(\Sigma_2)$ ; we have  $\text{char}(\Sigma_2) = \mathcal{M}_1$ .

## Generating Explanations

In this section, we consider the generation of all explanations for an atom  $q$ . We exclude in our considerations the trivial explanation  $E = \{q\}$ , which always exists if  $q$  belongs to the assumption literals  $A$ ,  $\Sigma \cup \{q\}$  is satisfiable and  $\Sigma \not\models q$ . These conditions can be efficiently checked under both formula- and model-based representations.

Recall that a prime implicate (res., prime implicant) of a theory  $\Sigma$  is a smallest (w.r.t. inclusion) clause  $c$  (resp., term  $t$ ) such that  $\Sigma \models c$  (resp.,  $t \models \Sigma$ ). As well-known, explanations can be characterized by prime implicates as follows.

**Proposition 2** *For a theory  $\Sigma$ ,  $E$  is a nontrivial explanation of  $q$  w.r.t.  $A \subseteq \text{Lit}$  if and only if the clause  $c = \bigvee_{p \in E} \bar{p} \vee q$  is a prime implicate of  $\Sigma$  such that  $E \subseteq A$ .*

We start with the generation of all nontrivial explanations under formula-based representation. For this problem, we present the following algorithm.

### Algorithm EXPLANATIONS

**Input:** A Horn CNF  $\varphi$  and a positive letter  $q$ .

**Output:** All nontrivial explanations of  $q$  from  $\varphi$ .

**Step 1.**  $\varphi^* := \emptyset$ ,  $S := \emptyset$ , and  $O := \emptyset$ ;

**Step 2. for each  $c \in \varphi$  do**

add any prime implicate  $c' \subseteq c$  of  $\varphi$  to  $\varphi^*$ ;  
**for each**  $c' \in \varphi^*$  with  $P(c') = \{q\}$  and  $N(c') \notin S$  **do**  
**begin** output  $N(c')$ ;  $S := S \cup \{N(c')\}$ ;  
 $O := O \cup \{(c, c') \mid c \in \varphi^*\}$   
**end;**

**Step 3. while** some  $(c_1, c_2) \in O$  exists **do**

**begin**  $O := O \setminus \{(c_1, c_2)\}$ ;  
**if** (1)  $q \notin N(c_1)$ , (2)  $P(c_1) = \{r\} \subseteq N(c_2)$  and  
(3)  $\varphi^* \cup N(c_1) \cup N(c_2) \setminus P(c_1)$  is satisfiable  
**then begin**  $c :=$  resolvent of  $c_1$  and  $c_2$ ;  
compute any prime implicate  $c' \subseteq c$  of  $\varphi$ ;  
**if**  $N(c') \notin S$  **then**  
**begin** output  $N(c')$ ;  $S := S \cup \{N(c')\}$ ;  
 $O := O \cup \{(c, c') \mid c \in \varphi^*\}$   
**end**  
**end**  
**end.** □

**Example 3** We consider algorithm EXPLANATION on input of  $\varphi = (\bar{x}_1 \vee \bar{x}_4)(\bar{x}_4 \vee \bar{x}_3)(\bar{x}_1 \vee x_2)(\bar{x}_3 \vee \bar{x}_5 \vee x_1)$ , and  $q = x_2$ . As easily seen, each clause in  $\varphi$  is a prime implicate, and thus after Step 2,  $\varphi^* = \varphi$  and  $S = \{\{x_1\}\}$ . Furthermore, the explanation  $E_1 = \{x_1\}$  was output.

In Step 3,  $(c_1, c_2)$ , where  $c_1 = \bar{x}_3 \vee \bar{x}_5 \vee x_1$  and  $c_2 = \bar{x}_1 \vee x_2$ , is the only pair in  $O$  satisfying (2)  $P(c_2) = \{x_1\} \subseteq N(c_2) (= \{x_1\})$ ; moreover, (1)  $q \notin N(c_1) (= \{x_3, x_5\})$  holds and (3)  $\varphi^* \cup \{x_3, x_5\}$  is satisfiable. Thus, in the body of the while-loop, its resolvent  $c = \bar{x}_3 \vee \bar{x}_5 \vee x_2$  is computed.

Clause  $c$  is a prime implicate of  $\varphi$ , and thus  $E_2 = \{x_3, x_5\}$  is output and added to  $S$ . Furthermore,  $O$  is updated.

In the next iterations, no pair  $(c_1, c_2) \in O$  is found which satisfies condition (2), and thus the algorithm halts. Note that  $E_1$  and  $E_2$  are the nontrivial explanations of  $q = x_2$ . □

The following result states that our algorithm works as desired. For any formula  $\varphi$ , denote by  $\|\varphi\|$  its length, i.e., the number of literal occurrences in it.

**Theorem 1** *Algorithm EXPLANATIONS incrementally outputs, without duplicates, all nontrivial explanations of  $q$  from  $\varphi$ . Moreover, the next output resp. termination occurs within  $O(e \cdot m \cdot n \cdot \|\varphi\|)$  time, where  $m$  is the number of clauses in  $\varphi$ ,  $n$  the number of atoms, and  $e$  the number of explanations output so far.*

**Proof.** (Sketch) Only pairs  $(c, c')$  are added to  $O$  such that  $c'$  is a prime implicate of  $\varphi$ . Furthermore, by condition (3) in Step 3, each such  $c'$  must have  $P(c') = \{q\}$ . Thus, by Props. 1 and 2, algorithm EXPLANATIONS outputs only nontrivial (clearly different) explanations  $E_1, E_2, \dots, E_k$  for  $q$ .

To show that it outputs all nontrivial explanations  $E$  for  $q$ , assume that such an  $E$  is not output, i.e.,  $E \neq E_i, i = 1, \dots, k$ . Let  $\varphi_N$  be the CNF of all negative prime implicates of  $\varphi$ , and let  $\varphi' = \varphi_N \cup \{\bigvee_{p \in E_i} \bar{p} \vee q \mid i \in \{1, \dots, k\}\}$ . Since  $c = \bigvee_{p \in E} \bar{p} \vee q$  is a prime implicate of  $\varphi$  and  $c \notin \varphi'$ , there exists a model  $v \in \text{mod}(\varphi')$  such that  $c(v) = 0$ . Let  $w$  be a maximal such model, i.e., no model  $u (> w)$  exists such that  $u \in \text{mod}(\varphi')$  and  $c(u) = 0$ . Since  $c(w) = 0$  implies  $w \notin \text{mod}(\varphi)$ , there exists a prime implicate  $c_1$  in  $\varphi^*$  (when Step 2 is finished) such that  $c_1(w) = 0$ . Clearly  $c_1 \notin \varphi'$ , i.e.,  $c_1$  is of form  $c_1 = \bigvee_{p \in N(c_1)} \bar{p} \vee x_i$  such that  $x_i \neq q$ . Moreover, we have  $q \notin N(c_1)$  by  $c(w) = c_1(w) = 0$ . Consider now the model  $w'$  defined by  $w'_i = 1$  and  $w'_j = w_j$ , for all  $j \neq i$ . Note that  $c(w') = 0$ , and by the maximality of  $w$ , there is a prime implicate  $c_2 \in \varphi'$  such that  $c_2(w') = 0$ . Since  $c_2(w) = 1$  and  $c_2(w') = 0$ , we have  $x_i \in N(c_2)$ . Since  $q \notin N(c_1)$ , a resolvent  $c^*$  of  $c_1$  and  $c_2$  thus exists. It can be shown (Eiter & Makino 2002) that  $c^*$  creates a new prime implicate  $c' = \bigvee_{p \in N(c')} \bar{p} \vee q \subseteq c^*$  of  $\varphi$  in Step 3, i.e.,  $N(c') \neq E_i, i = 1, \dots, k$ . This contradicts our assumption.

Thus, EXPLANATION is correct, and it remains to verify the time bound. Computing a prime implicate  $c' \subseteq c$  of  $\varphi$  in Steps 2 and 3 is feasible in time  $O(n \cdot \|\varphi\|)$ , and thus the outputs in Step 2 occur with  $O(m \cdot n \cdot \|\varphi\|)$  delay. As for Step 3, note the  $O$  contains only pairs  $(c_1, c_2)$  where  $c_1 \in \varphi^*$  and  $c_2 = N(c_2) \cup \{q\}$  such that  $N(c_2)$  was output, and each such pairs is added to  $O$  only once. Thus, the next output or termination follows within  $e \cdot m$  runs of the while-loop, where  $e$  is the number of solutions output so far. The body of the loop can be done, using proper data structures, in  $O(n \cdot \|\varphi\|)$  time (for checking  $N(c_1) \notin S$  efficiently, we may store  $S$  in a prefix tree). Thus, the time until the next output resp. termination is bounded by  $O(e \cdot m \cdot n \cdot \|\varphi\|)$ . □

From this result, we obtain the following corollary.

**Corollary 1** *Given a Horn CNF  $\varphi$  and a query  $q$ , computing  $O(n^k)$  many explanations of  $q$ , where  $k$  is a constant, is possible in polynomial time.*

This corollary implies that Selman and Leveque’s conjecture (1996, p. 266) that generating  $O(n)$  many explanations of  $q$  is NP-hard, where  $n$  is the number of propositional letters in the language, is not true (unless  $P=NP$ ). Note, however, that by the results of (Selman & Levesque 1996), computing  $O(n)$  many or all assumption-based explanations from a Horn  $\Sigma$  is not possible in total polynomial time unless  $P = NP$ .

Let us now consider computing all explanations in the model-based setting.

**Theorem 2** *Given the characteristic set  $\mathcal{M} = \text{char}(\Sigma) \subseteq \{0, 1\}^n$  of a Horn theory  $\Sigma$ , a query  $q$ , and  $A \subseteq \text{Lit}$ , computing the set of all explanations for  $q$  from  $\Sigma$  w.r.t.  $A$  is polynomial-time equivalent to dualizing a positive CNF.*

Here, polynomial-time equivalence means mutual polynomial-time transformability between deterministic functions, i.e.,  $A$  reduces to  $B$ , if there a polynomial functions  $f, g$  s.t. for any input  $I$  of  $A$ ,  $f(I)$  is an input of  $B$ , and if  $O$  is the output for  $f(I)$ , then  $g(O)$  is the output of  $I$ , cf. (Papadimitriou 1994); we also request that  $O$  has size polynomial in the size of the output for  $I$  (if not, trivial reductions may exist). In our reduction, explanations correspond to clauses of the dual prime CNF and vice versa.

**Proof.** (Sketch) By Props. 1 and 2, we need only compute all nontrivial explanations  $E \subseteq A$  corresponding to prime implicates  $c$  of  $\Sigma$  s.t.  $P(c) = \{q\}$  and  $N(c) = E \subseteq A \cap P$ .

We describe how the problem can be transformed to dualization of polynomially many positive CNFs  $\varphi_1, \dots, \varphi_k$ , such that the clauses of the dual prime CNFs  $\psi_i$  for  $\varphi_i$  correspond to the explanations of  $q$  w.r.t.  $A \cap P$  (equivalently, w.r.t.  $A$ ). Thus, the problem is polynomially reducible to dualizing (in parallel) several positive CNFs  $\varphi_i$ . By simple methods, we can combine  $\varphi_1, \dots, \varphi_k$  into a single CNF  $\varphi$  (using further variables) such that the clauses of the dual prime CNF for  $\varphi$  correspond to the explanations of  $q$  w.r.t.  $A$ . (This step is of less interest in practice, since dualization of the individual  $\varphi_i$  is at the core of the computation.)

To construct the  $\varphi_i$ , we proceed as follows. Let  $q = x_j$ .

(1) Define  $\mathcal{M}_i = \{v \in \mathcal{M} \mid v_j = i\}$ ,  $i \in \{0, 1\}$ .

(2) For every model  $v \in \max(\mathcal{M}_1)$ , let

$$F_v = \max(\{v \wedge x^A \wedge w \mid w \in \max(\mathcal{M}_0)\}).$$

We associate with  $F_v$  a monotone Boolean function  $f_v$  on the variables  $P_v = A \cap \{x_i \mid v_i = 1\}$  such that  $f_v(w) = 0 \Leftrightarrow w \leq s$  for some vector  $s$  in the projection of  $F_v$  on  $P_v$ . That is,  $F_v$  describes the maximal false points of  $f_v$ .

(3) Finally, define for every  $v \in \max(\mathcal{M}_1)$

$$\varphi_v = \{c \mid N(c) = \emptyset, P(c) = P_v \setminus S, x^S \in F_v\}.$$

Note that  $\varphi_v$  is a prime CNF for  $f_v$ .

It can be shown (Eiter & Makino 2002) that the nontrivial explanations of  $q$  w.r.t.  $A$  are given by the clauses in all dual prime CNFs  $\psi_v$  for  $\varphi_v$  where  $v \in \max(\mathcal{M}_1)$  (equivalently, by all prime implicants  $t$ , i.e., a prime DNF representation of  $f_v$ ,  $v \in \max(\mathcal{M}_1)$ ). This proves one direction of the result.

For the converse, we show that, given a positive CNF  $\varphi$  on atoms  $P$ , computing an equivalent prime DNF  $\psi$  is reducible

to computing all explanations as follows. Let  $q$  be a fresh letter (for component  $n + 1$ ), and define  $\mathcal{M} = \{(v, 0) \mid v \in \max(\{w \mid \varphi(w) = 0\})\} \cup \{(11 \dots 1)\}$  and  $A = P$ ; note that  $\max(\{w \mid \varphi(w) = 0\})$  is easily computed from  $\varphi$ .  $\square$

**Example 4** Let  $\mathcal{M} = \{(11011), (11010), (10101), (01010), (00001)\}$ , and suppose we want all explanations of  $q = x_1$  w.r.t.  $A = \{x_3, x_4, x_5\}$ . According to above, we obtain:

(1)  $\mathcal{M}_0 = \{(01010), (00001)\}$  and  $\mathcal{M}_1 = \{(11011), (11010), (10101)\}$ , thus  $\max(\mathcal{M}_1) = \{(11011), (10101)\}$ .

(2) We have two vectors  $v^{(1)} = (11011)$  and  $v^{(2)} = (10101)$ :

$$\begin{aligned} F_{v^{(1)}} &= \max \left( \left\{ \begin{array}{l} (11011) \wedge (00111) \wedge (01010), \\ (11011) \wedge (00111) \wedge (00001) \end{array} \right\} \right) \\ &= \{(00010), (00001)\}, \\ F_{v^{(2)}} &= \max \left( \left\{ \begin{array}{l} (10101) \wedge (00111) \wedge (01010), \\ (10101) \wedge (00111) \wedge (00001) \end{array} \right\} \right) \\ &= \{(00001)\}. \end{aligned}$$

Thus,  $P_{v^{(1)}} = \{x_4, x_5\}$  and  $f_{v^{(1)}}(w) = 0$  iff  $w \in \{(10), (01), (00)\}$ , and  $P_{v^{(2)}} = \{x_3, x_5\}$  and  $f_{v^{(2)}}(w) = 0$  iff  $w \in \{(01), (00)\}$ .

(3) We obtain  $\varphi^{(1)} = x_4 \wedge x_5$  and  $\varphi^{(2)} = x_3$ . The respective prime dual CNFs are  $\psi^{(1)} = x_4 \vee x_5$  and  $\psi^{(2)} = x_3$ .

Thus, the explanations of  $q$  w.r.t.  $A$  are  $E_1 = \{x_4, x_5\}$  and  $E_2 = \{x_3\}$ . It can be seen that this is the correct result.

## Negative Queries

So far, we considered Horn theories  $\Sigma$  and queries given by a letter  $q$ . In a general setting, we might allow that the formulas in  $\Sigma$  and the query are any propositional formulas. As for computation, we can introduce for a query, given by any formula  $\chi$ , a fresh letter  $q$ , add implications  $q \rightarrow \chi$ ,  $\chi \rightarrow q$  in  $\Sigma$ , and then ask for a nontrivial explanation of  $q$ . Thus, positive letter queries do not constrain the expressivity of the framework. However, this technique does not work for Horn theories, if one of the implications  $q \rightarrow \chi$ ,  $\chi \rightarrow q$  is not Horn. In the simplest case,  $\chi$  is a negative literal  $\bar{q}$ .

The next result tells us that already in this case, abduction from a Horn CNF is intractable. Recall that a Horn CNF  $\varphi$  is *acyclic*, if the graph on  $P$  with arcs from  $x_i \in N(c)$  to  $x_i \in P(c)$ ,  $c \in \varphi$ , has no directed cycle.

**Theorem 3** *Given a Horn CNF  $\varphi$ , a general query  $\chi$  in CNF, and  $A \subseteq \text{Lit}$ , deciding if  $\chi$  has a nontrivial explanation w.r.t.  $A$  is NP-complete. Hardness holds even if  $\chi = \bar{q}$ ,  $\varphi$  is acyclic, and either (i)  $A = P$  or (ii)  $A = P' \cup \bar{P}'$  for some  $P' \subseteq P$ .*

**Proof.** (Sketch) The problem is in NP, since clearly an explanation  $E$  exists if some set  $E \subseteq A$  exists such that  $\Sigma \cup E$  is satisfiable and  $\Sigma \cup E \models \chi$ ; such an  $E$  can be guessed and the conditions can be checked in polynomial time.

Hardness is shown by a reduction from 3SAT. Let  $\gamma = c_1 \wedge \dots \wedge c_m$  be a 3CNF over atoms  $x_1, \dots, x_n$ , where  $c_i = \ell_{i,1} \vee \ell_{i,2} \vee \ell_{i,3}$ . We introduce for each clause  $c_i$  a new atom  $y_i$ , for each  $x_j$  a new atom  $x'_j$ , and a special atom  $z$ . The Horn CNF  $\varphi$  contains the following clauses:

- $\bar{x}_i \vee \bar{x}'_i$ , for all  $i = 1, \dots, n$ ;
- $\bar{z} \vee y_1$
- $\bar{y}_i \vee \bar{\ell}_{i,j} \vee y_{i+1}$  if  $\ell_{i,j}$  is positive and  $\bar{y}_i \vee \bar{\ell}'_{i,j} \vee y_{i+1}$  if  $\ell_{i,j}$  is negative, for all  $i = 1, \dots, m-1$  and  $j = 1, 2, 3$ ;
- $\bar{y}_m \vee \bar{\ell}_{m,j}$  if  $\ell_{m,j}$  is positive and  $\bar{y}_m \vee \bar{\ell}'_{m,j}$  if  $\ell_{i,j}$  is negative, for  $j = 1, 2, 3$ .

As easily seen,  $\varphi$  is acyclic Horn. It can be shown (Eiter & Makino 2002) that the query  $q = \neg z$  has a nontrivial explanation  $E$  consisting of positive literals iff  $\gamma$  is satisfiable, which proves NP-hardness under restriction (i). For (ii), we use a similar construction.  $\square$

Note that this result contrasts the tractability result that a nontrivial explanation  $E \subseteq P$  for a positive query  $q$  can be computed in polynomial time (Selman & Levesque 1996). Thus, the framework of Horn abduction is sensitive with respect to query representation. We also remark that we can find an arbitrary explanation  $E$  for a query  $\bar{q}$  (which may contain negative literals), in polynomial time.

In the model-based setting, we obtain for computing all explanations for a negative query a similar result as for a positive query.

**Theorem 4** *Given the characteristic set  $\mathcal{M} = \text{char}(\Sigma) \subseteq \{0, 1\}^n$  of a Horn theory  $\Sigma$ , a negative query  $\bar{q}$ , and  $A \subseteq \text{Lit}$ , computing all explanations for  $q$  from  $\Sigma$  w.r.t.  $A$  is polynomial-time equivalent to dualizing a positive CNF.*

**Proof.** (Sketch) Observe that since the query is not positive, explanations of  $q$  may involve negative literals.

Proposition 2 implies that the nontrivial explanations for  $\bar{q}$  w.r.t.  $A$  correspond to the prime implicates  $c$  of  $\Sigma$  such that  $q \in N(c)$  and  $P(c) \cup N(c) \subseteq A \cup \{q\}$ . Let  $q = x_j$ , and define sets  $\mathcal{M}_0$  and  $\mathcal{M}_1$  for  $\mathcal{M}$  as in the proof of Theorem 2. Denote by  $A_+$  (resp.,  $A_-$ ) the set of positive (resp., negative) literals in  $A$ . We consider the following two cases:

(1) *Positive explanations for  $\bar{q}$ .* I.e., all prime implicates  $c$  of  $\Sigma$  s.t.  $\{q\} \subseteq N(c) \subseteq A_+ \cup \{q\}$  and  $P(c) = \emptyset$ .

Similarly as in the proof of Theorem 2, we construct dualization problems for functions  $f_v$ , but for  $v \in \max(\mathcal{M}_0)$ :

(1.1) For every  $v \in \max(\mathcal{M}_0)$ , let

$$F_v = \max(\{v \wedge x^{A_+} \wedge w \mid w \in \max(\mathcal{M}_1)\}).$$

The associated monotone Boolean function  $f_v$  on  $P_v = A \cap \{x_i \mid v_i = 1\}$  is defined by  $f_v(w) = 0 \Leftrightarrow w \leq s$  holds for some vector  $s$  in the projection of  $F_v$  on  $P_v$ .

(1.2) We define, for  $v \in \max(\mathcal{M}_0)$ ,

$$\varphi_v = \{c \mid N(c) = \emptyset, P(c) = P_v \setminus S, x^S \in F_v\}.$$

Similarly as in Theorem 2, we can show that the clauses in the dual prime CNFs for all  $\varphi_v, v \in \max(\mathcal{M}_0)$ , correspond to the positive explanations of  $\bar{q}$  (Eiter & Makino 2002).

(2) *Non-positive explanations for  $\bar{q}$ .* These are all prime implicates  $c$  of  $\Sigma$  s.t.  $\{q\} \subseteq N(c) \subseteq A_+ \cup \{q\}$  and  $P(c) = \{r\}$ , where  $r \in A_-$ .

For each  $r = x_{j'}$  (where  $j' \neq j$ ), we proceed as follows.

(2.1) For every  $v \in \max(\mathcal{M}_0)$  and  $i \in \{0, 1\}$ , define

$$\mathcal{M}_i^r = \{v \in \mathcal{M}_i \mid v_r = i\}.$$

(2.2) For each  $v^{(0)} \in \max(\mathcal{M}_0^r)$  and  $v^{(1)} \in \max(\mathcal{M}_1^r)$ , let

$$F_{v^{(0)}, v^{(1)}} = \max(\{v^{(0)} \wedge v^{(1)} \wedge x^{A_+} \wedge w \mid w \in \max(\{u \in \mathcal{M}_1 \mid u_{j'} = 0\})\}).$$

We associate with it a monotone Boolean function  $f_{v^{(0)}, v^{(1)}}$  on  $P_{v^{(0)}, v^{(1)}} = A \cap \{x_i \mid v_i^{(0)} = v_i^{(1)} = 1\}$  such that  $f_{v^{(0)}, v^{(1)}}(w) = 0 \Leftrightarrow w \leq s$  for some  $s$  in the projection of  $F_{v^{(0)}, v^{(1)}}$  on  $P_{v^{(0)}, v^{(1)}}$ .

(2.3) We define the CNFs

$$\varphi_{v^{(0)}, v^{(1)}} = \{c \mid N(c) = \emptyset, P(c) = P_v \setminus S, x^S \in F_{v^{(0)}, v^{(1)}}\}.$$

Then, it can be shown (Eiter & Makino 2002) that the clauses in the dual prime CNFs  $\psi_{v^{(0)}, v^{(1)}}$  for all  $\varphi_{v^{(0)}, v^{(1)}}$ , where  $v^{(0)} \in \max(\mathcal{M}_0^r)$  and  $v^{(1)} \in \max(\mathcal{M}_1^r)$  and  $r \in A_-$ , correspond to the non-positive explanations of  $\bar{q}$ .

In total, computing all explanations of  $\bar{q}$  is polynomial-time reducible to dualizing (in parallel) polynomially many positive CNFs. As mentioned in the proof of Theorem 2, this is polynomially reducible to dualizing a single CNF.

The converse is shown by a reduction similar to the one in the proof of Theorem 2; we just invert the polarity of  $q$ .  $\square$

## Joint Explanations

We call any set of  $E \subseteq A$  of literals a *joint explanation* of observations  $o_1, o_2, \dots, o_l$  from a background theory  $\Sigma$  w.r.t. a set of assumptions  $A \subseteq \text{Lit}$ , if  $E$  is an explanation of each  $o_i$  from  $\Sigma$  w.r.t.  $A$ . The observations  $o_i$  may be letters, or in a generalized setting propositional formulas.

Note that any such  $E$  is also an explanation for the conjunction  $\alpha = o_1 \wedge o_2 \wedge \dots \wedge o_l$  of all observations, while the converse is not true in general: an explanation  $E$  of  $\alpha$  may not satisfy minimality for  $o_1$ , say, i.e., some  $E' \subset E$  may explain  $o_1$ . Thus, joint explanations are stronger than ordinary explanations. In case of multiple explanations, this may be used to single out those which match with each of the (possibly independently made) observations.

For example, the malfunctioning of a car may be explained by two car mechanics, based on observations  $o_1$  and  $o_2$ , respectively. A match of their (individual) diagnoses  $E_1$  and  $E_2$  (i.e.,  $E_1 = E_2$ ) may be taken in favor of believing in their correctness. In fact, the diagnoses are robust in the sense that adding the other observation does not require a change; from another perspective, the same diagnosis is good for explaining different observations. If  $E_1$  and  $E_2$  are different, then we might want to know whether alternative diagnoses  $E'_1$  and  $E'_2$  do exist which coincide, i.e., whether a joint explanation is possible.

As it turns out, recognizing joint explanations for CNFs, i.e., deciding whether  $E$  is a joint explanation for observations  $o_1, \dots, o_l$  described by CNFs, from  $\Sigma$  w.r.t. assumptions  $A$  is tractable, for both formula- and model based representation. However, deciding existence is harder.

**Theorem 5** *Given a Horn CNF  $\varphi$ , query CNFs  $\chi_1, \chi_2, \dots, \chi_l$ , where  $l \geq 2$ , and  $A \subseteq \text{Lit}$ , deciding if a joint explanation exists from  $\Sigma$  w.r.t.  $A$  is NP-complete. Hardness holds even if  $l = 2$ , each  $\chi_i$  is a letter,  $\varphi$  is acyclic, and  $A = \text{Lit}$ .*

**Theorem 6** Given the characteristic set  $\mathcal{M} = \text{char}(\Sigma)$  of a Horn theory  $\Sigma$ , query CNFs  $\chi_1, \chi_2, \dots, \chi_l$ ,  $l \geq 2$ , and  $A \subseteq \text{Lit}$ , deciding if a joint explanation exists from  $\Sigma$  w.r.t.  $A$  is NP-complete. Hardness holds even if  $l = 2$ , each  $\chi_i$  is a letter, and  $A = \text{Lit}$ .

Thus, the tractability results in (Selman & Levesque 1996; Kautz, Kearns, & Selman 1993) do not generalize to joint explanations for positive queries. Similar intractability results hold for negative queries and combined positive and negative queries.

### Related Works

Selman and Levesque (1990; 1996) were among the first to study the complexity of computing general and assumption-based explanations; Corollary 1 closes an open problem of them. The underlying algorithm EXPLANATIONS is a relative of a similar procedure by Boros *et al.* (1990) for computing all prime implicates of a Horn CNF in output-polynomial time. In fact, Theorem 1 can be seen as a strengthening of their result. For negative queries, a similar algorithm is not evident. del Val (2000) presented generation of implicates and prime implicates of certain clausal theories in a target language, which is formed on a subset of the atoms, using a procedure based on kernel resolution and derived exponential bounds on its running time. Furthermore, del Val described the use of this procedure for generating jointly all explanations of all literals not on a set of atoms  $V$ . However, neither is this method incremental in nature, nor is it clear whether it is total polynomial time. Moreover, it considers a letter  $q$  and its negation  $\bar{q}$  at once.

Inoue (1992) considered, in the propositional and the first-order context, generating explanations and prime implicates using SOL-resolution. He proposed a strategy which processes, starting from the empty set, clauses from a theory incrementally. However, due to possible large intermediate results, this method is not total polynomial time in general. Khardon *et al.* (1999) show how computing all keys of a relational database schema, which is constrained by a Boolean formula  $\varphi$ , can be polynomially transformed into computing all explanations of a query  $q$  from  $\varphi \wedge \psi$ , where  $\psi$  is Horn. Thus, our algorithm EXPLANATIONS can be used for efficiently generating all keys of a database scheme where  $\varphi$  is a Horn CNF. This generalizes the result for  $\varphi$  consisting of non-negative Horn clauses, i.e., a set of functional dependencies (Lucchesi & Osborn 1978).<sup>1</sup> Note that Khardon *et al.* also show how to compute a single explanation of a query  $q$  from a theory  $\varphi$  polynomially, using repeatedly an oracle for computing a key of a database schema constrained by  $\varphi \wedge \psi$ , where  $\psi$  is Horn; however, this method is not usable for generating explanations from general Horn theories (cf. Footnote 1). Less related to our work is (Eiter & Gottlob 1995), which considered abduction from Horn and general propositional theories, but focused on existence of explanations and reasoning tasks about explanations.

### Conclusion

We have presented a number of positive and negative results about generating all and some abductive explanations, re-

<sup>1</sup>In fact, Khardon *et al.*'s transformation works only if  $\varphi$  has no negative prime implicate; otherwise, keys introduce inconsistency.

spectively, which complement previous work in the literature. In particular, we analyzed the role of positive vs negative abductive queries, under both formula- and model-based representation, and we considered the novel notion of joint explanation. Our positive results may be readily applied for efficiently computing (a subset of) all explanations. The results draw a complete picture for the model-based setting, and almost so for the formula-based setting; the complexity of generating all explanations for a negative query in it is currently open.

### Acknowledgments

We thank the reviewers for helpful comments. This work was partly supported by the Austrian Science Fund (FWF) project Z29-INF, by TU Wien through a scientific collaboration grant, and by the Scientific Grant in Aid of the Ministry of Education, Science, Sports and Culture of Japan.

### References

- Bioch, C., and Ibaraki, T. 1995. Complexity of identification and dualization of positive Boolean functions. *Information and Computation* 123:50–63.
- Boros, E.; Crama, Y.; and Hammer, P. L. 1990. Polynomial-time inference of all valid implications for Horn and related formulae. *Annals of Mathematics & Artificial Intelligence* 1:21–32.
- Brewka, G.; Dix, J.; and Konolige, K. 1997. *Nonmonotonic Reasoning – An Overview*. Number 73 in CSLI Lecture Notes. CSLI Publications, Stanford University.
- Dechter, R., and Pearl, J. 1992. Structure identification in relational data. *Artificial Intelligence* 58:237–270.
- del Val, A. 2000. The complexity of restricted consequence finding and abduction. In *Proc. AAAI-00*, 337–342.
- Eiter, T., and Gottlob, G. 1995. The complexity of logic-based abduction. *J. ACM* 42(1):3–42.
- Eiter, T., and Makino, K. 2002. On computing all abductive explanations (preliminary report). Technical Report INFVSYS RR-1843-02-04, Institut für Informationssysteme, TU Wien, Austria.
- Eiter, T.; Gottlob, G.; and Makino, K. 2002. New results on monotone dualization and generating hypergraph transversals. In *Proc. 34th ACM Symp. on Theory of Computing (STOC 2002)*.
- Fredman, M., and Khachiyan, L. 1996. On the complexity of dualization of monotone disjunctive normal forms. *Journal of Algorithms* 21:618–628.
- Inoue, K. 1992. Linear resolution for consequence finding. *Artificial Intelligence* 56(2-3):301–354.
- Kautz, H.; Kearns, M.; and Selman, B. 1993. Reasoning with characteristic models. In *Proc. AAAI-93*, 34–39.
- Khardon, R., and Roth, D. 1996. Reasoning with models. *Artificial Intelligence* 87(1/2):187–213.
- Khardon, R.; Mannila, H.; and Roth, D. 1999. Reasoning with examples: Propositional formulae and database dependencies. *Acta Informatica* 36(4):267–286.
- Lucchesi, C. L., and Osborn, S. 1978. Candidate keys for relations. *J. Computer and System Sciences* 17:270–279.
- Papadimitriou, C. H. 1994. *Computational Complexity*.
- Poole, D. 1989. Explanation and prediction: An architecture for default and abductive reasoning. *Comp. Intelligence* 5(1):97–110.
- Selman, B., and Levesque, H. J. 1990. Abductive and default reasoning: A computational core. In *Proc. AAAI-90*, 343–348.
- Selman, B., and Levesque, H. J. 1996. Support set selection for abductive and default reasoning. *Artif. Intell.* 82:259–272.