

# A Superposition Decision Procedure for the Guarded Fragment with Equality

Harald Ganzinger\*

Hans de Nivelle\*\*

Max-Planck-Institut für Informatik, D-66123 Saarbrücken

{hg | nivelle}@mpi-sb.mpg.de

## Abstract

*We give a new decision procedure for the guarded fragment with equality. The procedure is based on resolution with superposition. We argue that this method will be more useful in practice than methods based on the enumeration of certain finite structures. It is surprising to see that one does not need any sophisticated simplification and redundancy elimination method to make superposition terminate on the class of clauses that is obtained from the clausification of guarded formulas. Yet the decision procedure obtained is optimal with regard to time complexity. We also show that the method can be extended to the loosely guarded fragment with equality.*

## 1 Introduction

The loosely guarded fragment was introduced in (Andréka, van Benthem & Németi 1996) as ‘the modal fragment of classical logic’. It is obtained essentially by restricting quantification to the following forms:

$$\forall y[R(x, y) \rightarrow A(x, y)] \text{ and } \exists y[R(x, y) \wedge A(x, y)].$$

These forms naturally arise when modal formulae are translated into classical logic using the standard translation based on the Kripke frames. The authors showed there that the guarded fragment has many of the nice properties of modal logics. In particular it is decidable. Any decision procedure for this fragment, hence, is a decision procedure for those modal logics that can be embedded into it, for example  $K$ ,  $D$ ,  $S3$ , and  $B$ . It has been shown by Grädel (1997) that equality can be admitted in the guarded fragment without affecting decidability. In the fragment with equality additional logics such as difference logic can be expressed (where  $\diamond A$  means  $A$  holds in a world different from the present).

De Nivelle (1998) has given a resolution decision procedure for the guarded fragment without equality. In his procedure, a non-liftable ordering is employed, and, hence, some additional and non-trivial argument is required for

proving refutational completeness. In this paper we describe in detail a decision procedure for the guarded fragment with equality which is based on resolution and superposition. Despite the fact that it applies to a larger fragment of first-order logic, our new procedure is simpler than the one in (de Nivelle 1998) in that we employ a liftable ordering (plus selection) so that we are able to re-use standard results about refutational completeness. Our method is also interesting as there are not so many saturation-based decision procedures for fragments with equality described in the literature. Notable exceptions include (Fermüller & Salzer 1993), where a resolution decision procedure is given for the Ackermann class with equality, and (Bachmair, Ganzinger & Waldmann 1993), where it is shown that a certain superposition strategy decides the monadic class with equality. Nieuwenhuis (1996) proves the decidability of certain shallow equational theories by basic paramodulation.

The advantage of resolution or superposition decision procedures over theoretical procedures based on collapsing models is that the former use syntactic, unification-based inferences to enumerate candidate witnesses of inconsistency. There is experimental evidence (Hustadt & Schmidt 1997) that such procedures perform well in practice, in particular they often will not exhibit the usually exponential or double-exponential worst-case complexity of the respective fragments. Also, when having a flexible saturation theorem prover at hand, such as SPASS (Weidenbach 1997), it suffices to appropriately adjust its parameters in order to efficiently implement the procedure.

The results of this paper can be summarized as follows. (i) Ordered paramodulation with selection is a decision procedure for the GF with equality. No sophisticated redundancy elimination methods are required, and a straightforward (liftable) ordering and selection strategy suffice. (ii) The procedure decides the class of guarded clauses which is a proper superclass of the GF with equality. (iii) The worst-case time complexity of the decision procedure is doubly-exponential, which is optimal, given that the logic is 2EXPTIME-complete (Grädel 1997). (iv) Guarded clauses with deep terms, although decidable in the case without equality, become undecidable in the equational case. (v) The superposition-based decision method can be extended

\*Work supported in part by the ESPRIT Basic Research Working Group 22457 (CCL II).

\*\*Work done at ILLC, U. Amsterdam, Plantage Muidergracht 24, 1018 TV Amsterdam

to the loosely guarded fragment with equality, but is much more involved there. For the extension, hyper-inferences which simultaneously resolve a set of guards are needed. Some non-trivial results are required about the existence of suitable partial inferences to avoid the generation of clauses which are not loosely guarded, together with meta-theorems about the refutational completeness of these partial inferences.

## 2 The Guarded Fragment

**Definition 2.1** The formulas of the *guarded fragment*  $\mathbf{GF}$  of *function-free* first order logic are inductively defined as follows:

1.  $\top$  and  $\perp$  are in  $\mathbf{GF}$ .
2. If  $A$  is an atom then  $A$  is in  $\mathbf{GF}$ .
3.  $\mathbf{GF}$  is closed under boolean combinations.
4. If  $F \in \mathbf{GF}$  and  $G$  is an atom, for which every free variable of  $F$  is among the arguments of  $G$ , then  $\forall \bar{x}(G \rightarrow F) \in \mathbf{GF}$  (or, equivalently,  $\forall \bar{x}(\neg G \vee F) \in \mathbf{GF}$ ) and  $\exists \bar{x}(G \wedge F) \in \mathbf{GF}$ , for every sequence  $\bar{x}$  of variables.

The atoms  $G$  which appear as constraints for quantified variables are called *guards*. Equations can also be used as guards. These are examples of guarded formulae:

$$\begin{aligned} & \forall x (x \approx x \rightarrow p(x)), \quad \exists x (p(x) \wedge q(x)) \\ & \forall yz (r(y, y, z) \rightarrow \perp), \quad \forall xy (r(x, y) \rightarrow r(y, x)) \\ & \quad \forall xy (r(x, y) \rightarrow \exists z r(y, z)) \\ & \exists x [R(w, x) \wedge \forall y (R(x, y) \rightarrow p(y)) \wedge q(x)] \end{aligned}$$

The last formula is the translation of the modal formula  $\diamond(\Box p \wedge q)$  with respect to a world  $w$ . These are formulae which are not guarded:

$$\begin{aligned} & \forall xy p(x, y) \\ & \forall x_1 x_2 x_3 [p(x_1, x_2) \rightarrow p(x_2, x_3) \rightarrow p(x_1, x_3)]. \end{aligned}$$

The last formula states the transitivity of  $p$ . As this is not guarded, for modal logics such as  $S4$  which are based on transitive frames the standard embedding methods lead outside the guarded fragment.

## 3 The Superposition Calculus

For the decision procedure to be described below we only need a rather weak form of the superposition calculus of Bachmair & Ganzinger (1990), called ordered paramodulation, for which Hsiang & Rusinowitch (1991) have also given a completeness proof. Here (ordered) paramodulation into the larger side of an equation is permitted. We use the symbol  $\approx$  to denote formal equality and do not distinguish between equations  $s \approx t$  and  $t \approx s$ . Disequations  $\neg(s \approx t)$  will also be written as  $s \not\approx t$ . The calculus is clausal, where clauses are multisets of literals  $L_1, \dots, L_k$ ,  $k \geq 0$ , which we write as disjunctions  $L_1 \vee \dots \vee L_k$ . A clause is called *positive* if it does not contain any negative literals. A clause is called *ground* or *propositional* if it contains no variables.

The calculus is parameterized by admissible orderings  $\succ$  and selection functions  $\Sigma$  for negative literals. For each setting of the two parameters it is refutationally complete. For dealing with the orderings it is useful to view non-equational atoms of the form  $p(t_1, \dots, t_k)$ , with  $p$  a predicate symbol different from equality, as a shorthand notation for an equation  $p(t_1, \dots, t_k) \approx \text{tt}$ . In this encoding, the atom is considered a term (in a two-sorted signature with sorts  $i$  and  $o$ ), with  $\text{tt}$  a distinguished constant of sort  $o$ , and where predicates are viewed as functions of sort  $o$ , taking arguments of sort  $i$ . An *admissible ordering*  $\succ$  is any total reduction ordering on ground terms (including non-equational ground atoms) in which  $\text{tt}$  is minimal. The multiset extension of  $\succ$ , again denoted  $\succ$ , is used to compare literals by identifying any positive equation  $s \approx t$  (including the equational encodings of non-equational atoms) with the multiset  $\{s, t, \text{tt}\}$ , and any negative equation  $s \not\approx t$  with the multiset  $\{s, t, \text{tt}\}$ , respectively. The ordering is extended to non-ground expressions by defining  $E \succ E'$  iff, for all ground substitutions  $\sigma$ ,  $E\sigma \succ E'\sigma$ . Although admissible orderings are total and well-founded on ground terms and literals, they are only partial on non-ground expressions. Whenever a literal  $L$  contains a unique maximal term we will denote it by  $\max(L)$ . A *selection function*  $\Sigma$  selects, in each clause, at most one (occurrence of a) negative literal. This occurrence is called *selected*.

Inferences involve eligible literals. A literal is called *eligible* in a clause  $C$  if either it is selected in  $C$  (by  $\Sigma$ ), or else nothing is selected in  $C$ , and it is a maximal literal in  $C$  with respect to  $\succ$ . In particular, a positive literal, since it cannot be selected, is eligible only if the respective clause contains no selected (negative) literal. The inference rules are as follows:

**Ordered Factoring.** From  $A_1 \vee A_2 \vee R$  derive  $A_1\sigma \vee R\sigma$  provided  $A_1$  is eligible and  $\sigma$  is the mgu of  $A_1$  and  $A_2$ .

**Equality Factoring.** From  $t_1 \approx u \vee t_2 \approx v \vee R$  derive  $u\sigma \not\approx v\sigma \vee t_1\sigma \approx v\sigma \vee R\sigma$  provided  $t_1 \approx u$  is eligible and  $\sigma$  is the mgu of  $t_1$  and  $t_2$ .

**Reflexivity Resolution.** From  $t_1 \not\approx t_2 \vee R$  derive  $R\sigma$  provided that  $t_1 \not\approx t_2$  is eligible and  $\sigma$  is the mgu of  $t_1$  and  $t_2$ .

**Resolution.** From  $A_1 \vee R_1$  and  $\neg A_2 \vee R_2$  derive  $R_1\sigma \vee R_2\sigma$  provided that both  $A_1$  and  $\neg A_2$  are eligible and  $\sigma$  is the mgu of  $A_1$  and  $A_2$ .

**Ordered Paramodulation.** From  $t_1 \approx u \vee R_1$  and  $L[t_2] \vee R_2$ , where  $t_2$  is not a variable, derive  $L[u]\sigma \vee R_1\sigma \vee R_2\sigma$  provided that both  $t_1 \approx u$  and the literal  $L[t_2]$  are eligible,  $\sigma$  is the mgu of  $t_1$  and  $t_2$ , and  $u \not\approx t_1$ .

The way in which the order restrictions are applied here is *a priori*, i.e. before the unifier is computed. Superposition is complete also if the order restrictions are checked after the

substitution is applied to the premises (*a posteriori* checking), or even if they are attached to the clauses and inherited throughout inferences. A priori checking has the advantage that the eligible literals in a clause can be precomputed, before any inference is attempted. On the other hand, a posteriori application is generally more restrictive. For obtaining the theoretical results in the present paper a priori ordering constraints turn out to be sufficiently powerful.

The calculus is refutationally complete for any choice of admissible ordering and selection function. Moreover, the calculus is compatible with a rather powerful notion of redundancy by which don't-care non-deterministic simplification and redundancy elimination can be justified. In particular, tautologies can be eliminated and multiple occurrences of literals in clauses can be deleted. The notion of redundancy allows for much more sophisticated simplification methods which, however, will not be required here, although for achieving good practical performance they have to be implemented. The fact that non-naive implementations of superposition, such as in the SPASS system, spend most of their execution time on simplification rather than search is what makes them useful in the end. We call a set of clauses *N saturated up to redundancy* (with respect to ordered paramodulation) if any inference from non-redundant premises in *N* is redundant in *N*. The definition of redundancy, in particular, implies that an inference is redundant in *N* if the conclusion of the inference is contained in *N* or else is redundant in *N*.

**Theorem 3.1 (Bachmair & Ganzinger, 1990)** Let *N* be a set of clauses that is saturated up to redundancy with respect to the above derivation rules. Then *N* is unsatisfiable if and only if *N* contains the empty clause.

## 4 The Decision Procedure

We will now describe the decision procedure. We define a notion of guarded clauses, and show that guarded formulae can be translated into guarded clause sets. We will obtain a resolution decision procedure by defining a reduction order  $\succ$  and a selection function  $\Sigma$  that force an upper bound on the complexity of the derivable clauses.

### 4.1 Clausal Normal Form Translation

We rely on a specific clausal normal-form transformation for the guarded fragment. We may assume that the given formula is in negation normal form, that is, negation is only applied to atoms. We also assume that implications and equivalences have been eliminated by replacing them by equivalent formulas involving conjunction, disjunction, and negation. These standard transformations do not take a formula outside the guarded fragment.

The next step is to replace certain sub-formulae by fresh names, together with a definition of the name.<sup>1</sup> We abstract

<sup>1</sup>Such transformations are called structural and are, for instance, studied in (Baaz, Fermüller & Leitsch 1994). They are called structural

universally quantified sub-formulae to reduce the number of quantifier alternations. Let  $\mathcal{F} = \{F_1, \dots, F_n\}$  be a set of formulae in negation normal form. The *structural transformation* of  $\mathcal{GF}$  is obtained by iterating the following transformations: If *F* is a formula in  $\mathcal{F}$  containing a proper sub-formula of the form  $\forall \bar{x}(-G \vee H)$ , with *G* the guard, then (i) add a *definition*  $\forall \bar{x}\bar{y}(-G \vee \neg \alpha(\bar{y}) \vee H)$  to  $\mathcal{F}$ , and (ii) replace the indicated sub-formula in *F* by  $\alpha(\bar{y})$ . Hereby it is assumed that  $\bar{y}$  is the set of variables that occur in *G*, but not in  $\bar{x}$ , and that  $\alpha$  is a new predicate name that does not occur in  $\mathcal{F}$ . Observe that the structural transformation, when applied to a set of guarded formulas also yields a set of guarded formulas as result. Moreover, all remaining universal quantifiers are outermost, so that any inner existential quantifier occurs in the scope of all universally quantified variables. This method of eliminating embedded quantifiers is standard and has also been used in the context of the guarded fragment by Grädel (1997).

For the purposes of this paper, the standard skolemization technique is the one which is appropriate. One replaces any applied occurrence of an existentially quantified variable *y* by a term  $f(x_1, \dots, x_n)$ , with *f* a new Skolem function symbol, if  $x_1, \dots, x_n$  are the universally quantified variables, in the scope of which *y* occurs. After that replacement, all existential quantifiers have been removed, and Skolem function applications contain all the variables of a formula. Finally, to obtain a set of clauses, distribute disjunctions over conjunctions, omit the universal quantifiers (which are all outermost) and consider any conjunction of disjunctions as a set (of clauses).

**Example 4.1** Consider the guarded formula

$$\exists x (n(x) \wedge \forall y [\neg a(x, y) \vee \forall z \{\neg p(x, z) \vee \exists x (a(x, z) \wedge (\neg b(z, z) \vee \neg c(x, x))\})]).$$

The structural transformation gives the set of formulas

$$\begin{aligned} &\exists x [n(x) \wedge \alpha(x)], \\ &\forall x, y [\neg a(x, y) \vee \neg \alpha(x) \vee \beta(x)], \\ &\forall x, z [\neg p(x, z) \vee \neg \beta(x) \vee \\ &\quad \exists x (a(x, z) \wedge (\neg b(z, z) \vee \neg c(x, x)))]. \end{aligned}$$

Skolemization yields

$$\begin{aligned} &n(c) \wedge \alpha(c), \\ &\forall x, y [\neg a(x, y) \vee \neg \alpha(x) \vee \beta(x)], \\ &\forall x, z [\neg p(x, z) \vee \neg \beta(x) \vee \\ &\quad (a(fxz, z) \wedge (\neg b(z, z) \vee \neg c(fxz, fxz)))]. \end{aligned}$$

Clauserification, finally, produces this set of clauses:

$$\begin{aligned} &n(c) \\ &\alpha(c) \\ &\neg a(x, y) \vee \neg \alpha(x) \vee \beta(x) \\ &\neg p(x, z) \vee \neg \beta(x) \vee a(fxz, z) \\ &\neg p(x, z) \vee \neg \beta(x) \vee \neg b(z, z) \vee \neg c(fxz, fxz). \end{aligned}$$

since more of the structure of a formula is preserved when the formula is factored.

## 4.2 Guarded Clauses

The result of the transformation are sets of guarded clauses which, in particular, consist of a specific kind of literals. A term is called *shallow* if either it is a variable or else a functional term  $f(u_1, \dots, u_m)$ ,  $m \geq 0$ , in which each  $u_j$  is a variable or a constant. A literal  $L$  is called *simple* if each term in  $L$  is shallow. Hence  $p(x, c, f(x))$  and  $f(x, c) \not\approx y$  are simple while  $\neg p(s(f(0), x))$  and  $f(x, s(x)) \approx g(x)$  are not. A clause is called *simple* if all literals are simple. A literal is called *covering* if each non-ground and non-variable subterm in the literal contains all the variables of the literal. An expression is called *functional* if it contains a constant or a function symbol, and *non-functional*, otherwise.

**Definition 4.2** A simple clause  $C$  is called *guarded* if it satisfies the following conditions:

(i)  $C$  is a positive, non-functional, single-variable clause; or

(ii) every functional subterm in  $C$  contains all the variables of  $C$ , and, if  $C$  is non-ground,  $C$  contains a non-functional negative literal, called a *guard*, which contains all the variables of  $C$ .

Clauses of the form (ii) are called *properly guarded*, while the concept of guards is void for the other types of guarded clauses. A set of clauses is called *guarded* if all its clauses are guarded.

Note that if a guarded clause contains a constant it must be a ground clause in which terms are shallow. Also, any literal in a guarded clause is covering.

These are some examples of guarded clauses where suitable guards have been underlined.

$$\begin{aligned} & p(0, s(0)) \vee c \not\approx d \vee q(s(0), f(0, 0)) \\ & \underline{p(x, x)} \vee q(x) \\ & \neg p(y, x) \vee \underline{\neg q(x, y, y)} \vee r(x + y, x - y, x) \\ & \underline{\neg p(y, x)} \vee \neg q(x, y, y) \\ & \underline{x \not\approx y} \vee x \approx (x + y) \end{aligned}$$

The following clauses are not guarded:

$$\begin{aligned} & \neg e(x) \vee e(s(x)) && \text{(not simple)} \\ & \neg p(x) \vee \neg q(y) \vee r(x, y) && \text{(no guard)} \\ & \neg p(f(x, y)) \vee p(x, y) && \text{(no guard)} \\ & \neg p(x, y) \vee p(f(x), y) && \text{(not covering)} \\ & \neg p(x, y) \vee p(0, g(x, y)) && \text{(constant, but non-ground)} \end{aligned}$$

Definition 4.2 is more restrictive than the corresponding definition in (de Nivelles 1998). The last two clauses in the previous example are guarded in the sense of (de Nivelles 1998). In the section 5 we will discuss this issue in more detail.

**Theorem 4.3** The number of different (up to variable renaming) guarded clauses (without duplicate occurrences of literals) over a finite signature has a double exponential upper bound in the size of the signature.

*Proof.* Let a finite signature be given. Define the following parameters:

$a_1$	the maximal arity of function symbols
$a_2$	the maximal arity of predicate symbols
$a$	the maximum of $a_1$ and $a_2$
$n_1$	$a_2 +$ the number of constant and function symbols
$n_2$	the number of predicate symbols
$n$	the maximum of $n_1$ and $n_2$ .

The maximal size  $s$  of a simple atom is  $a^2 + a + 1$ . Therefore, the number of simple atoms (modulo variable renaming) that may appear in a guarded clause over the given signature is bounded by

$$n^s = n^{a^2+a+1}.$$

Then the number of simple literals (modulo variable renaming) is at most

$$l = 2n^{a^2+a+1}.$$

This is also an upper bound for the maximal number of literals in a clause, since a clause contains at most all possible literals over at most  $a_2$  variables. Then the number of guarded clauses that can be constructed from non-repeated literals is bounded by

$$c = 2^l = 2^{2n^{a^2+a+1}}.$$

□

## 4.3 Preservation of Guardedness

We now show that guarded clauses are closed under the paramodulation inferences so that, using the theorem 4.3, saturating a given set of clauses under these inferences, combined with eager elimination of duplicate literals in clauses, yields a decision procedure for satisfiability. To that end we need to define an appropriate ordering and selection function. For the ordering  $\succ$  we may use any lexicographic path ordering on terms and non-equational atoms based on a precedence  $\succ$  such that  $f \succ c \succ p \succ \text{tt}$  for any non-constant function symbol  $f$ , constant  $c$ , and predicate symbol  $p$ , respectively. For the selection function  $\Sigma$  we assume that (i) if a clause is non-functional and contains a guard then one of its guards is selected by  $\Sigma$ ; (ii) if a clause contains a functional negative literal, one of these is selected; and (iii) if a clause contains a positive functional literal, but no negative functional literal, no literal is selected, so that the maximality principle applies for a literal to be eligible for an inference.

**Lemma 4.4** Let  $L_1, L_2$  be two literals of a guarded clause. Assume that  $L_2$  contains a non-ground functional term, while  $L_1$  does not. Then  $L_2 \succ L_1$ .

*Proof.* First observe that with the given assumptions the clause does not contain any constants. Let  $L_1$  be a literal, and let  $t$  be a functional term in  $L_2$ . First suppose that  $L_1$  is a non-equational literal of the form  $[\neg]p(u_1, \dots, u_n)$  with

variables  $u_i$ . Then, any of the  $u_i$  also occurs in  $t$ . With regard to the ordering, non-equational literals such as  $L_1$  are identified with equations  $[\neg](p(u_1, \dots, u_n) \approx \text{tt})$ . Let  $f$  be the leading function symbol in  $t$ . Then  $f$  has a precedence greater than any of the symbols in  $L_2$ , and as  $t$  contains all variables of  $L_1$ , we conclude that  $p(u_1, \dots, u_n) < t \preceq \max(L_2)$  which implies that  $L_1 < L_2$ .

If  $L_1$  is an equational atom  $u \approx v$ , by a similar reasoning we infer that  $t > u$  and  $t > v$ , from which again  $L_1 < L_2$  is inferred.  $\square$

**Lemma 4.5** With  $>$  and  $\Sigma$  as defined above, a literal in a clause is eligible for an inference only if it contains all the variables of the clause.

**Lemma 4.6** Let  $\sigma$  be the most general unifier of two simple non-equational atoms  $p(t_1, \dots, t_n)$  and  $p(u_1, \dots, u_n)$ . Then  $p(t_1, \dots, t_n)\sigma$  is also simple.

**Lemma 4.7** Let  $A$  and  $B$  be simple atoms such that (i) every variable occurring in  $B$  also occurs in  $A$ ; (ii) every variable that occurs in a functional term of  $B$  also occurs in a functional term of  $A$ ; and (iii) every functional term of  $B$  contains all the variables of  $A$ . Then for any substitution  $\sigma$ ,

- (i) if  $A\sigma$  is simple, then  $B\sigma$  is simple,
- (ii) every variable of  $B\sigma$  occurs in  $A\sigma$ ,
- (iii) every variable occurring in a functional term of  $B\sigma$  occurs in a functional term of  $A\sigma$ .
- (iv) Every functional term of  $B\sigma$  contains all the variables of  $A\sigma$ .

As a consequence of the lemma 4.4, if a clause is non-ground, any eligible literal either contains a (non-ground) functional term or else there is no functional term in the entire clause. The preceding lemma can therefore be applied to any eligible literal  $A$  and any other literal  $B$  in a guarded clause.

**Lemma 4.8** A factor of a guarded clause is guarded.

**Lemma 4.9** An equality factor of a guarded clause is guarded.

**Lemma 4.10** A clause obtained by reflexivity resolution from a guarded clause, is guarded.

*Proof.* The propositional case the lemma is trivial. For reflexivity resolution to be applicable to a non-propositional clause, the clause must be of the form  $D = x \not\approx y \vee C$ , with guard  $x \not\approx y$  and with  $C$  not containing a functional term. Clearly, the resolvent has only simple literals and is either the empty clause or has just one variable. In the latter case the resolvent either has a guard or is a positive clause.  $\square$

**Lemma 4.11** A resolvent of two guarded clauses is guarded.

*Proof.* Let  $C_1 = A_1 \vee D_1$  and  $C_2 = \neg A_2 \vee D_2$  be the clauses resolved upon, with  $\sigma$  the mgu of  $A_1$  and  $A_2$ . Then the conclusion is the clause  $D = D_1\sigma \vee D_2\sigma$ . Notice that with  $A = A_i$  and  $B$  any literal in  $D_i$ , the premises of the lemma 4.7 are satisfied, both for  $i = 1$  and  $i = 2$ . As both  $A_1$  and  $A_2$  are simple, the literal  $A_1\sigma$  is also simple. Applying the lemma 4.7, part (i), we may infer that all literals in  $D$  are simple. If there are functional terms in  $D$  then these contain the same set of variables, and all the variables of  $D$ , cf. Theorem 4.7, parts (iii) and (iv). In order to show that there is a guard in  $D$  when one is needed, we distinguish as to whether or not the clauses are ground.

Suppose that one of the  $C_i$  is ground. In that case  $D$  is ground since literals which are eligible for an inference contain all the variables of a clause.

Let us now assume that both  $C_1$  and  $C_2$  are non-ground. Suppose that  $C_1$  is not a positive clause over one variable. Then  $C_1$  must have a guard  $\neg G$ , and  $\neg G\sigma$  occurs in  $D$ . Moreover,  $A_1$  must have a functional term containing all the variables of  $C_1$ . (Otherwise  $\neg G$  or some other guard of  $C_1$  would be selected and the inference would not be possible). As  $A_1\sigma$  is simple,  $\sigma$  assigns a variable to each variable in  $C_1$ . Therefore, the literal  $\neg G\sigma$  has only variables as arguments. Since  $\neg G\sigma$  contains all the variables of  $A_1\sigma$ , it contains all variables of  $D$ , and, hence, is a guard. In case that  $C_1$  is a positive, single-variable clause, then  $D$  contains at most one variable. If there is no guard in  $D$  then the resolvent must be a single-variable, positive, possibly empty clause.

Finally, the resolvent does not contain a constant unless one of the premises does. In that case both the premise and the resolvent are ground.  $\square$

**Lemma 4.12** Any clause obtained by a superposition inference from two guarded clauses is guarded.

*Proof.* Let  $C_1 = L[u] \vee D_1$  be the main premise,  $C_2 = t_1 \approx t_2 \vee D_2$  the side premise, and  $D = L[t_2]\sigma \vee D_1\sigma \vee D_2\sigma$  be the conclusion, respectively, of the inference, with  $\sigma$  the mgu of  $t_1$  and  $u$ .

We first consider the case where  $C_2$  is ground. If  $t_2$  is not a constant then also  $t_1$  is not a constant, as otherwise the ordering constraints would block the inference. Superposition inferences into variables are excluded so that  $u$  must be a functional term containing all the variables of the clause. Hence, all variables in  $u$  become grounded by  $\sigma$ ,  $D$  is ground, and contains simple literals only.

If  $C_2$  is non-ground, then  $t_1 \approx t_2$  has to contain all its variables, and at least one of the  $t_1$  or  $t_2$  is a functional term. (Otherwise the guard in  $C_2$  would be selected and the clause cannot appear as the side premise of the inference.) The ordering restrictions, therefore, imply that  $t_1$  is functional, containing all the variables of the clause, whereas  $t_2$  can be

a variable, or a functional term. The possible forms of  $u$  are also restricted.  $u$  cannot be a variable.  $u$  can be a functional term containing all the variables of  $C_2$ , or a ground term. Suppose that  $u$  is ground and unifiable with  $t_1$ . Then  $u$  is not a constant,  $C_2$  is ground, and  $u$  occurs as an argument to a predicate in  $C_2$ . Then,  $D$  is a ground clause and is simple since  $t_2\sigma$  is either a constant or a functional term with constant arguments. If  $u$  is not ground  $\sigma$  is a variable renaming and, in particular, both  $D_1\sigma$  and  $D_2\sigma$  are guarded. Moreover,  $L[t_2]\sigma$  is simple. It is easily checked that the guards of  $C_1\sigma$  and  $C_2\sigma$  can both serve as guards of  $D$ .  $\square$

**Theorem 4.13** Let  $\Sigma$  and  $\succ$  be as specified. For all the inferences of the ordered paramodulation calculus, if the premises are guarded, so is the conclusion.

**Theorem 4.14** The fragment of guarded clauses is decidable by ordered paramodulation.

*Proof.* By the theorem 4.13 all derivable clauses are guarded, and the number of such clauses is finite, cf. Theorem 4.3. As each inference rule is a decidable relation on guarded clauses, the theorem follows.<sup>2</sup>  $\square$

The theorem can also be extended to guarded clauses combined with unrestricted ground clauses. There one replaces in the initial clause set any ground (sub-) term  $s$  which is not shallow by a new constant  $a_s$ , together with the defining equation  $a_s \approx s$ . This preserves satisfiability and produces a clause set which is guarded.

#### 4.4 Complexity

The complexity of our decision procedure is double exponential. Grädel (1997) has shown that the decision problem for the guarded fragment with equality is 2EXPTIME-complete, hence our procedure is theoretically optimal. We use the fact, cf. Theorem 4.3, that the number of guarded clauses has a doubly exponential bound and show that the saturation process has no primitive operation that has more than exponential complexity.

**Theorem 4.15** The superposition decision procedure can be implemented in 2EXPTIME (in the size of the signature).

*Proof.* We reuse the notation defined in the proof of the theorem 4.3. It is clear that the space complexity of the procedure is dominated by the space that is needed to store the clauses. Hence, we obtain a space complexity of  $s * l * c$ .

<sup>2</sup>The inferences are equipped with constraints which specify which literals are eligible for an inference. Depending on the signature, the term ordering, and the selection function such constraints are in general undecidable and have to be approximated. This is not the case here. But even if the constraints were undecidable, by Theorem 4.13 a safe approximation would be to consider any unrestricted inference the conclusion of which is a guarded clause.

For the time complexity, observe that suitable abstractions of the ordering and selection constraints for the inferences can be checked in polynomial time, cf. the proof of Theorem 4.14. Then one may show that the time needed to do a subsumption check is in  $O(l^3s)$ . In fact, one first matches the guard with at most  $l$  literals. After that one has to try to match each of the  $l$  remaining literals with one of the  $l$  literals of the other clause. This gives a total of  $l^3$  attempted matches. Since each matching can take up to  $s$  time, this number has to be multiplied by  $s$ . Knowing the time complexity for subsumption for guarded clauses, we can estimate the time complexity of our method as a whole. The algorithm has to try all pairs of literals, and in the case that a resolvent is possible, it has to check that the resolvent is not subsumed by one of the existing clauses. This takes time in  $O((cl)^2c(l^3s))$ . This iteration has to be repeated at most  $c$  times, resulting in a bound in  $O((cl)^2c^2l^3s)$ . This number is roughly equal to  $c^4$  which gives the desired double exponential time complexity.

Finally we should also consider the time and space complexity of the clausal normal form translation. It is well-known that the transformation to normal form can take at most single exponential time, which is negligible compared to the double exponential time obtained above. The (structural) elimination of equivalences is slightly more tricky here as the result has to be a guarded formula.  $\square$

## 5 Weakly Guarded Clauses

The notion of guarded clause as given in the Definition 4.2 is more restrictive than the one given in (de Nivelle 1998). There, terms of arbitrary depth are allowed provided that they are either ground, or contain all variables of the clause. We repeat the formal definition:

**Definition 5.1** A clause  $C$  is called *weakly guarded*, if (i) every non-ground functional term in  $C$  contains all the variables of  $C$ ; and (ii) if  $C$  is non-ground it contains a negative literal, all of which arguments are constants or variables, and which contains all the variables of the clause.

This notion was inspired by the  $E^+$ -class. Every clause which is guarded is also weakly guarded, but the converse is not true in general.

**Theorem 5.2** Satisfiability is undecidable for finite sets of weakly guarded clauses if equational atoms are admitted. The fragment remains undecidable if all ground terms are constants.

The Post Correspondence Problem can be reduced to this decision problem. This is essentially due to the fact that projection functions defined by equations of the form  $f(x, y) \approx x$  can make a non-shallow term equal to a term that violates the covering condition. For example from the guarded clauses  $\neg p(x, y) \vee p(s(f(x, y)))$  and  $\neg p(x, y) \vee f(x, y) \approx x$  we may deduce the non-guarded

clause  $\neg p(x, y) \vee p(s(x))$ , where  $s$  is not applied to all the variables of the clause. This shows that variables in nested functional terms cannot be combined with equality.

## 6 The Loosely Guarded Fragment

Our method can be generalized to the so called *loosely guarded fragment*. This fragment obtained by weakening the condition (4) in the Definition 2.1 as follows: If  $F$  is loosely guarded and  $G_1, \dots, G_n$  are atoms, with variables as arguments, then the formulae  $\forall \bar{x}(G_1 \wedge \dots \wedge G_n \rightarrow F)$  and  $\exists \bar{x}(G_1 \wedge \dots \wedge G_n \wedge F)$  are loosely guarded, provided that (i) every free variable of  $F$  occurs in a  $G_i$ , and (ii) every pair of variables  $y_1, y_2$ , which are free in  $F$ , and of which at least one is among the  $\bar{x}$ , occur together in one of the  $G_i$ . We call the entire conjunction  $G_1 \wedge \dots \wedge G_n$  the *guard* of the formula, and any conjunct a *guard atom*.

In the loosely guarded fragment the until operator can be expressed, which cannot be expressed in the guarded fragment.  $P$  until  $Q$  can be translated as:

$$\exists y (Rxy \wedge Qy \wedge \forall z (Rxz \wedge Rzy \rightarrow Pz)).$$

Transitivity of  $R$ , though, cannot be expressed in the loosely guarded fragment. In the formula

$$\forall x, y, z (Rxy \wedge Ryz \rightarrow Rxz)$$

there is no atom in the guard in which the variables  $x$  and  $z$  co-occur. In fact, Ganzinger, Meyer & Veanes (1999) have shown that allowing for a single transitive relation makes the LGF undecidable in general.

A CNF transformation similar to the one described in the section 4.1 leads to what we call loosely guarded clauses:

**Definition 6.1** A simple clause  $C$  is called *loosely guarded* if it satisfies the following conditions:

- (i)  $C$  is a positive, non-functional, single-variable clause; or
- (ii)  $C$  contains no constants, every functional subterm in  $C$  contains all the variables of  $C$ , and  $C$  contains a set of negative, non-functional literals  $\neg A_1, \dots, \neg A_n$ ,  $n \geq 0$ , called a (*loose*) *guard* of  $C$ , such that every pair of variables that occurs in  $C$  occurs together in one of the atoms  $A_i$ .

Propositional simple clauses are admitted. They have an empty guard.

The main modification of the decision procedure is that in cases where previously a guard atom needed to be selected in a clause now a set of literals may constitute a guard, and some of these have to be resolved simultaneously. Therefore, resolution needs to be generalized to (ordered) hyper-resolution. The basis for this are more general selection functions  $\Sigma$  which now may select an entire, possibly empty set of occurrences of negative literals in a clause. Now a literal is called *selected* if it occurs in the set of selected literals of a clause.

## Ordered Hyper-Resolution with Selection

$$\frac{A_1 \vee R_1 \quad \dots \quad A_k \vee R_n \quad \neg B_1 \vee \dots \vee \neg B_n \vee R}{R_1 \sigma \vee \dots \vee R_n \sigma \vee R \sigma}$$

where (i) either the  $\neg B_j$  are the literals selected by  $\Sigma$  in the *main premise*, or else  $n = 1$ , nothing is selected in  $\neg B_1 \vee R$ , and  $\neg B_1$  is maximal in  $\neg B_1 \vee R$ , (ii) the  $A_i$  are eligible in the *side premises*  $A_i \vee R_i$ , and (iii)  $\sigma$  is the mgu of the tuples  $(A_1, \dots, A_n)$  and  $(B_1, \dots, B_n)$ .

Given a hyper-resolution inference of this form, we speak of a *partial inference* producing a *partial conclusion*  $D$  whenever there exists a non-empty subset  $j_1, \dots, j_k$  of the indices  $1 \leq j \leq n$  and

$$D = \bigvee_{1 \leq i \leq k} R_{j_i} \tau \vee \bigvee_{i \notin \{j_1, \dots, j_k\}} \neg B_i \tau \vee R \tau,$$

with  $\tau$  the mgu of  $(A_{j_1}, \dots, A_{j_k})$  and  $(B_{j_1}, \dots, B_{j_k})$ .

The extended calculus is refutationally complete and compatible with a notion of redundancy by which the usual simplification mechanisms (tautology elimination, condensation, subsumption) can be justified. There is no published result that exactly covers this calculus, but it is easy to generalize the results in (Bachmair & Ganzinger 1990) appropriately.

The orderings which we may use for the decision procedure are the same as for the non-loose case. The selection function  $\Sigma$  should satisfy these restrictions:

- (i) If a clause  $C$  is non-functional and contains a guard  $L_1 \vee \dots \vee L_k$  then *all* the literals of one of the guards of  $C$  are selected by  $\Sigma$ ;
- (ii) if a clause contains a functional negative literal, *one* of these is selected; and
- (iii) if a clause contains a positive functional literal but no negative functional literal, then no literal is selected, so that the maximality principle applies for a literal to be eligible for an inference.

In order to prove that with this ordering and selection strategy, ordered paramodulation becomes a decision procedure for the LGF, two problems have to be solved. The first problem is that conclusions of inferences might become too deep.

**Example 6.2 (de Nivelle & de Rijke, 1999)** The following clause  $D$  is loosely guarded:

$$\neg a_1(x, y) \vee \neg a_2(y, z) \vee \neg a_3(z, x) \vee b_1(x, y) \vee b_2(y, z) \vee b_3(z, x)$$

There are no functional terms, therefore the three guard literals are selected. The following three clauses are candidates for a hyperresolution inference:

$$\begin{aligned} C_1 &= \neg p_1(u) \vee a_1(fu, fu), \\ C_2 &= \neg p_2(v) \vee a_2(v, gv), \\ C_3 &= \neg p_3(w) \vee a_3(gw, w), \end{aligned}$$

From these one may derive the hyper-resolvent

$$\neg p_1(u) \vee \neg p_2(fu) \vee \neg p_3(fu) \vee \\ b_1(fu, fu) \vee b_2(fu, gfu) \vee b_3(gfu, fu),$$

with an mgu  $\sigma = [x, y, v, w := fu, z := gfu]$ . This resolvent has a non-shallow term which is not admitted for a loosely guarded clause.

A remedy to this problem is to resolve  $D$  only with a suitable subset of the side premises  $C_i$ . In the example, if we only resolve the second and third guard literal of  $D$  with  $C_2$  and  $C_3$ , respectively, we obtain the partial conclusion

$$\neg a_1(w, w) \vee \neg p_2(w) \vee \neg p_3(w) \\ \vee b_1(w, w) \vee b_2(w, gw) \vee b_3(gw, w).$$

The mgu of the partial inference is  $[y, v, x := w, z := gw]$ . This clause is loosely guarded, in particular, not too deep. It turns out that if an inference is possible then one of its partial conclusions will be a guarded clause. The proof makes use of the subsequent lemma which is a special case of a theorem in (de Nivelle & de Rijke 1999).

**Lemma 6.3** Let  $A_1, \dots, A_n$  and  $B_1, \dots, B_n$  be  $2n \geq 2$  simple literals such that

- (i) the  $B_i$  are non-functional;
- (ii) for all  $x, y$  in  $\text{Var}(B_1, \dots, B_n)$  there is a  $B_i$  such that  $x, y$  is in  $\text{Var}(B_i)$ ;
- (iii) the  $A_j$  are covering and functional;
- (iv)  $A_i$  and  $A_j$ , for  $i \neq j$ , have no common variables;
- (v) the  $A_i$  and the  $B_j$  have no common variables;
- (vi) the tuples  $(A_1, \dots, A_n)$  and  $(B_1, \dots, B_n)$  are unifiable.

Then there exists a non-empty subset  $j_1, \dots, j_k$  of the indices  $1 \leq j \leq n$  such that the tuples  $(A_{j_1}, \dots, A_{j_k})$  and  $(B_{j_1}, \dots, B_{j_k})$  are unifiable with an mgu  $\tau$  and

- (i) any of the  $A_{j_i} \tau$  ( $= B_{j_i} \tau$ ) is simple and covering;
- (ii) if  $x$  is a variable in any of the  $B_i$  or  $A_{j_i}$  and if  $y$  is a variable in  $y\tau$  then  $y$  also occurs in  $A_{j_i} \tau$ .

The proof which is given in full detail for the more general theorem in (de Nivelle & de Rijke 1999) is based on this observation: Let us assume, for simplicity, that the  $A_i$  are non-ground and that all non-constant symbols are binary. Then any of the  $A_i$  is of the form  $p(u, fuv)$ ,  $p(fuv, v)$ , or  $p(fuv, guv)$ , with more variants arising from exchanging  $u$  and  $v$  in one of the arguments of  $p$ ,  $f$ , and  $g$ . If we disregard the trivial one-variable case, any of the guard atoms is of the form  $p(x, y)$ , with different variables  $x$  and  $y$ . The problem of unifying all the  $A_i$  with the corresponding  $B_i$ , therefore, induces at least one unification problem of the form  $x = fxy$ ,  $y = fxy$ ,  $x = y$ ,  $fxy = fxy$ , or  $fxy = fyx$  on any pair  $x, y$  of variables in  $\text{Var}(B_1, \dots, B_n)$ . This is a consequence of the co-occurrence requirement (ii). If the unification problem is solvable with an mgu  $\sigma$  then if  $x$  is

in  $\text{Var}(B_1, \dots, B_n)$ , either  $x\sigma$  is of maximal depth  $d$  among all the  $y\sigma$ , for  $y$  in  $\text{Var}(B_1, \dots, B_n)$ , or else  $x\sigma$  is a subterm of some  $y\sigma$ , with  $y$  in  $\text{Var}(B_1, \dots, B_n)$ . Picking for the  $j_i$  those atoms in which a variable  $x$  with  $x\sigma$  of depth  $d$  appears, solves the problem. Any other variable in one of the  $B_{j_i}$  will be instantiated either by a term of the same depth and containing the same variables, or else by a direct subterm of a term of depth  $d$ .

The lemma covers exactly those unification problems which arise from hyper-resolution inferences with guard atoms  $B_i$  and corresponding positive atoms  $A_i$ . For the latter to be eligible for an inference they all have to contain a functional term. In other words, with the class of orderings  $\succ$  and selection functions  $\Sigma$  which we consider for the LGF, we obtain this theorem:

**Theorem 6.4** Suppose there is an inference by hyperresolution with respect to  $\succ$  and  $\Sigma$ . Then one of the partial inferences produces a (partial) conclusion which is a guarded clause.

The existence of suitable partial inferences solves our problem as the calculus remains complete if, for any potential hyper-inference from side premises  $C_1, \dots, C_k$  and main premise  $D$ , rather than deriving the full conclusion, we derive any don't-care non-deterministically chosen partial conclusion. A proof of this fact in the non-equational case has been given in (Bachmair & Ganzinger 1997), and the proof does not depend on any properties that are critical when adding equality. The criterion for which partial conclusion to choose is simply that the conclusion should be a guarded clause. With this modification of the calculus, the class of guarded clauses is closed under its inferences.

A second, simpler problem arises from the fact that loosely guarded clauses over any given finite signature may be arbitrarily long. Fortunately it is not difficult to see that the set of guarded clauses that can be derived with our inference system from an initially given finite set of guarded clauses is finite. This is an immediate consequence of the fact that the number of variables does not increase during an inference: The point here is that the loose guard of any generated clause is an instantiation of the loose guard of one of the parent clauses. Therefore, the number of variables in any derived clause is bound by the number of variables in one of the parent clauses.

**Lemma 6.5** If  $D$  is the [partial] conclusion of an inference from premises  $C_i$  then  $|\text{Var}(D)| \leq \max(|\text{Var}(C_i)|)$ .

Altogether we obtain:

**Theorem 6.6** Ordered Paramodulation with hyperresolution based on selection is a decision procedure for the LGF.

## 7 Conclusions

We have shown that it is possible to effectively decide the [loosely] guarded fragment with equality by superposition-based saturation provers. There is hope that usable decision procedures can be obtained from these results with existing standard theorem provers. This hope is supported by our theoretical optimality result (in the non-loose case) and by experimental evidence that has been obtained in using these theorem proving techniques in related application domains (Hustadt & Schmidt 1997). The GF has turned out to be a fragment of first-order logic with equality for which it is especially easy to configure superposition into an optimal decision procedure. Although the complexity issue has been neglected by and large in the literature on resolution-based decision procedures, we believe that in most cases of fragments which are complete for a particular time complexity class, the resolution-based methods can be implemented in this time bound. (Things are different for space complexity classes such as PSPACE and local theorem proving methods based on resolution and superposition where the reuse of space, as is standard with tableau methods, is not so straightforward.) The loosely guarded case is more tricky. However this paper also demonstrates that the theory of saturation-based theorem proving is sufficiently developed to be able to solve the problems without having to deal with technically difficult proof-theoretic arguments.

## References

- Andréka, H., van Benthem, J. & Németi, I. (1996), Modal languages and bounded fragments of predicate logic, Technical report, ILLC.
- Baaz, M., Fermüller, C. & Leitsch, A. (1994), A non-elementary speed up in proof length by structural clause form transformation, in 'Proc. LICS'94'.
- Bachmair, L. & Ganzinger, H. (1990), On restrictions of ordered paramodulation with simplification, in M. Stickel, ed., 'Proc. 10th Int. Conf. on Automated Deduction, Kaiserslautern', Vol. 449 of *Lecture Notes in Computer Science*, Springer-Verlag, Berlin, pp. 427–441.
- Bachmair, L. & Ganzinger, H. (1997), A theory of resolution, Research Report MPI-I-97-2-005, Max-Planck-Institut für Informatik, Saarbrücken, Saarbrücken. To appear in the *Handbook of Automated Reasoning*.
- Bachmair, L., Ganzinger, H. & Waldmann, U. (1993), Superposition with simplification as a decision procedure for the monadic class with equality, in G. Gottlob, A. Leitsch & D. Mundici, eds, 'Proc. of Third Kurt Gödel Colloquium, KGC'93', Vol. 713 of *Lecture Notes in Computer Science*, Springer-Verlag, Berlin, pp. 83–96. Revised version of Technical Report MPI-I-93-204.
- de Nivelle, H. (1998), A resolution decision procedure for the guarded fragment, in C. Kirchner & H. Kirchner, eds, 'CADE-15', pp. 191–204.
- de Nivelle, H. & de Rijke, M. (1999), A resolution decision procedure for the guarded fragment, Manuscript.
- Fermüller, C. G. & Salzer, G. (1993), Ordered paramodulation and resolution as decision procedure, in A. Voronkov, ed., 'Proceedings of the 4th International Conference on Logic Programming and Automated Reasoning (LPAR'93)', Vol. 698 of *LNAI*, Springer Verlag, St. Petersburg, Russia, pp. 122–133.
- Ganzinger, H., Meyer, C. & Veanes, M. (1999), The two-variable guarded fragment with transitive relations, in 'Proc. 14th IEEE Symposium on Logic in Computer Science', IEEE Computer Society Press, this volume.
- Grädel, E. (1997), On the restraining power of guards, Manuscript. To appear in the *Journal of Symbolic Logic*.
- Hsiang, J. & Rusinowitch, M. (1991), 'Proving refutational completeness of theorem proving strategies: The transfinite semantic tree method', *J. Association for Computing Machinery* **38**(3), 559–587.
- Hustadt, U. & Schmidt, R. A. (1997), On evaluating decision procedures for modal logics, in M. E. Pollack, ed., 'Proceedings of the Fifteenth International Joint Conference on Artificial Intelligence (IJCAI-97)', International Joint Conferences on Artificial Intelligence, Inc. (IJCAI) and Japanese Society for Artificial Intelligence (JSAI), Morgan Kaufmann, Nagoya, Japan, pp. 202–207.
- Nieuwenhuis, R. (1996), Basic paramodulation and decidable theories, in 'Proceedings of the Eleventh Annual IEEE Symposium On Logic In Computer Science (LICS'96)', IEEE Computer Society Press, pp. 473–483.
- Weidenbach, C. (1997), 'Spass version 0.49', *J. Automated Reasoning* **18**(2), 247–252.