

Terminating constraint set satisfiability and simplification algorithms for context- dependent overloading

**Rodrigo Ribeiro, Carlos Camarão &
Lucília Figueiredo**

**Journal of the Brazilian Computer
Society**

ISSN 0104-6500

Volume 19

Number 4

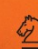
J Braz Comput Soc (2013) 19:423-432

DOI 10.1007/s13173-013-0107-9

Volume 19 · Number 4 · November 2013

Journal of the

► **Brazilian
Computer
Society** ►

 Springer

 Springer

Your article is protected by copyright and all rights are held exclusively by The Brazilian Computer Society. This e-offprint is for personal use only and shall not be self-archived in electronic repositories. If you wish to self-archive your article, please use the accepted manuscript version for posting on your own website. You may further deposit the accepted manuscript version in any repository, provided it is only made publicly available 12 months after official publication or later and provided acknowledgement is given to the original source of publication and a link is inserted to the published article on Springer's website. The link must be accompanied by the following text: "The final publication is available at link.springer.com".

Terminating constraint set satisfiability and simplification algorithms for context-dependent overloading

Rodrigo Ribeiro · Carlos Camarão ·
Lucília Figueiredo

Received: 24 August 2012 / Accepted: 4 March 2013 / Published online: 9 April 2013
© The Brazilian Computer Society 2013

Abstract Algorithms for constraint set satisfiability and simplification of Haskell type class constraints are used during type inference in order to allow the inference of more accurate types and to detect ambiguity. Unfortunately, both constraint set satisfiability and simplification are in general undecidable, and the use of these algorithms may cause non-termination of type inference. This paper presents algorithms for these problems that terminate on any given input, based on the use of a criterion that is tested on each recursive step. The use of this criterion eliminates the need of imposing syntactic conditions on Haskell type class and instance declarations in order to guarantee termination of type inference in the presence of multi-parameter type classes, and allows program compilation without the need of compiler flags for lifting such restrictions. Undecidability of the problems implies the existence of instances for which the algorithm incorrectly reports unsatisfiability, but we are not aware of any practical example where this occurs.

Keywords Haskell · Constraint set satisfiability · Constraint set simplification · Termination

1 Introduction

Haskell's type class system [5,18] extends the Hindley-Milner type system [16] with constrained polymorphic types, in order to support overloading. Type class constraints may occur in types of expressions involving overloaded names (or symbols), and restrict the set of types to which quantified type variables may be instantiated, to those types for which these type constraints are satisfied, according to types of definitions that exist in a relevant context.

A type class declaration specifies the name and parameters of the class, and the principal type of names which can then be overloaded in instance definitions. For example:

```
class Eq a where
  (==) :: a → a → Bool
  (/=) :: a → a → Bool
```

is a declaration of type class *Eq*, with parameter *a*, that specifies the principal types of `(==)` and `(/=)`. Function `(==)` has type $\forall a. Eq\ a \Rightarrow a \rightarrow a \rightarrow Bool$, where constraint *Eq a* indicates that type variable *a* cannot be instantiated to an arbitrary type, but only to a type that has been defined as an instance of class *Eq*.

An instance of a type class specifies instance types for type class parameters, and gives definitions of the overloaded names specified in the class. The type of each overloaded name in an instance definition is obtained by substituting type class parameters with corresponding instance types. For example, the following instance declarations specify definitions of the equality operator for types *Int* and for polymorphic lists, respectively:

```
instance Eq Int where
  (==) = primEqInt
```

R. Ribeiro (✉) · C. Camarão
Instituto de Ciências Exatas, Departamento de Ciência
da Computação, Universidade Federal de Minas Gerais,
Belo Horizonte, Brazil
e-mail: rodrigogribeiro@decea.ufop.br

C. Camarão
e-mail: camaraorribeiro@dcc.ufmg.br

L. Figueiredo
Instituto de Ciências Exatas e Biológicas,
Departamento de Computação,
Universidade Federal de Ouro Preto, Ouro Preto, Brazil
e-mail: lucilia@iceb.ufop.br

```

instance Eq a => Eq [a] where
  []      == []      = True
  (a:x) == (b:y) = a == b && x == y
  _      == _      = False

```

For a base type, like *Int*, a corresponding predefined operation is provided. The definition of equality for lists of elements of an arbitrary type uses the equality test for elements of this type. Constraint *Eq a* must be specified as the *context* for the *head* *Eq [a]* of the instance declaration. A *context* is a set of type class constraints, and constraint π is the *head* of a qualified constraint $P \Rightarrow \pi$, where P is a set of type class constraints.

As an aside, type classes in Haskell may also contain default definitions of the overloaded names, in order to avoid repeating the same definitions in instances.

Class constraints introduced on the types of overloaded symbols occur also on the types of expressions defined in terms of these symbols. For example, consider the following function that tests list membership:

```

elem a [] = False
elem a (b:x) = a == b || elem a x

```

The principal type of *elem* is $\forall a. Eq a \Rightarrow a \rightarrow [a] \rightarrow Bool$. Constraint *Eq a* occurs in the type of *elem* due to the use of the equality operator (*==*) in its definition.

Haskell restricts type classes to have a single parameter but the extension to multi-parameter type classes, called Haskell+mptcs in the sequel, is widely used.

Type inference for constrained type systems rely on constraint set simplification, which, for the case of type classes, essentially amounts to performing (so-called) *context reduction*. Constraint set simplification yields equivalent constraint sets, and are useful for providing simpler types for expressions. Context reduction simplifies constraints by substituting constraints or removing resolved constraints according to available instance definitions, besides removing duplicate constraints or substituting constraints according to the class hierarchy.

As an example, context *Eq [t]* is reduced to *Eq t*, for any type *t*, in the presence of instance *Eq [a]* with context *Eq a*.

Improvement [13] is also a process of simplification of constrained types, but it is of a different nature, and is used in type inference to avoid ambiguity and to infer more informative types. Improvement is fundamentally based on constraint set satisfiability: it is a process of transforming a constraint set P into a constraint set obtained by applying a substitution S to P so that the set of satisfiable instances of P is preserved.

The mechanism of functional dependencies and other alternatives have been proposed to deal with improvement [7,4,10,11,14], for detection of ambiguity and for specializa-

tion of constrained types in the presence of multi-parameter type classes. We do not discuss improvement specifically in this paper, but focus on constraint set satisfiability, which is only used for the implementation of improvement or any alternative approach.

Unfortunately, both constraint set satisfiability and simplification are in general undecidable problems [6], and the use of computable functions for solving these problems may cause non-termination of type inference.

This paper presents algorithms for constraint set satisfiability and simplification that use a termination criterion which is based on a measure of the sizes of types in type constraints. The sequence of constraints that unify with a constraint axiom in recursive calls of the function that checks satisfiability or simplification of a type constraint is such that either the sizes of types of each constraint in this sequence is decreasing or there exists at least one type parameter position with decreasing size.

The use of this criterion eliminates the need for imposing syntactic conditions on Haskell type class and instance declarations in order to guarantee termination of type inference in the presence of multi-parameter type classes, and allows program compilation without the need of compiler flags for lifting such restrictions.

The use of a termination criterion implies that there exist well-typed programs for which the presented algorithm incorrectly reports unsatisfiability. However, practical examples where this occurs are expected to be very rare. The algorithms have been implemented and tested by using a prototype front-end for Haskell, available at the mptc github repository. The algorithm works as expected when subjected to examples mentioned in the literature, Haskell libraries that use multi-parameter type classes and many tests, including those used by the most commonly used Haskell compiler [19], GHC, involving all pertinent GHC extensions.

Restrictions imposed on class and instance declarations in Haskell, in Haskell+mptcs and in GHC, and GHC compilation flags used to avoid these restrictions [20], are summarized in Sect. 2. Section 3 reviews entailment and satisfiability relations on type class constraints. Section 4 gives a definition of a computable function that returns the set of satisfiable substitutions of a given constraint set P , when it terminates. Subsection 4.1 defines a termination criterion and redefines this computable function in order to use this criterion. Section 5 defines a constraint set simplification computable function, based on the same termination criterion. Section 6 concludes.

2 Restrictions over class and instance declarations

This section summarizes the restrictions imposed on class and instance declarations in Haskell, Haskell+mptcs and in

GHC, and GHC compilation flags used to avoid these restrictions.

By default, GHC follows the Haskell language specification (i.e., the Haskell 98 report [8]), which imposes the following restrictions.

1. Each class declaration must have exactly one parameter.
2. The head of a qualified constraint in an instance declaration must have the form $C(T \bar{\alpha})$, where C denotes a class name, T a type constructor and $\bar{\alpha}$ a sequence of distinct type variables. Such overbar notation is used extensively in this paper: \bar{x} denotes a possibly empty sequence of elements in the set $\{x_1, \dots, x_n\}$, for some $n \geq 0$.
3. Each constraint in a context P of an instance declaration $P \Rightarrow C \bar{\tau}$ must have the form $C a$, where a is a type variable occurring in $\bar{\tau}$.

Restriction 1 allows only single-parameter type classes, but multi-parameter type classes are widely used by programmers and in Haskell libraries and are supported in many Haskell implementations. For example, consider type class *Map* parameterized by the key and element types, and the type class *Collection*, parameterized by the type constructor and the type of elements of the collection, partly sketched below:

```
class Eq a => Collection c a where
  empty :: c a
  insert, delete :: a -> c a -> c a
  member :: a -> c a -> Bool
  ...
```

`instance Show (Tree Int) where ...` is an example of an instance declaration that does not follow restriction (2), because the head of the constraint (which has an empty context) consists of type constructor *Tree* applied to *Int*, not to a type variable.

Flag `-XFlexibleInstances` can be used by GHC users to avoid enforcing condition (2), i.e., to allow the head of a constraint in an instance declaration to be arbitrarily nested. The next is an example that does not follow restriction (3), since $s a$ is not just a type variable: `instance Show (s a) => Show (Sized s a)...`

Instances that do not follow these restrictions are common in Haskell programs, specially in the presence of multi-parameter type classes.

Flag `-XFlexibleContexts` can be used by GHC users to avoid restriction (3). With the use of this flag, contexts are restricted as follows:

1. No type variable can have more occurrences in a constraint of a context than in the head.

2. The sum of the number of *occurrences* of type variables and type constructors in a context must be smaller than in the head.

This restriction is known as the *Paterson Condition*. In some cases, it is still over-restrictive. As an example, consider the following code:

```
data Rose f a = Rose (f (Rose f a))
```

```
instance (Show (f (Rose f a)), Show a) =>
  Show (Rose f a) where ...
```

This instance of *Show* is rejected by GHC because it has more occurrences of type variable f in a constraint than in the head. Flag `-XUndecidableInstances`, which lifts all restrictions (including those related to the use of functional dependencies), is needed to compile this code. With this flag, termination is ensured by imposing a depth limit on a recursion stack [20].

3 Constrained polymorphism and type class constraints

The Haskell type class system is based on the more general theory of *qualified types* [12], which extends the Hindley-Milner type system with constrained types.

The syntax of types with type class constraints is defined in Fig. 1, where meta-variable usage is also indicated. For simplicity, and following common practice, kinds are not considered explicitly in type expressions, and type applications are assumed to be well kinded. Function types $\tau_1 \rightarrow \tau_2$ are constructed as the curried application of the function type constructor to two arguments, and are written as usual in infix notation.

The union of constraint sets P and Q is denoted by P, Q and a slight abuse of notation is made by writing simply π for the singleton constraint set $\{\pi\}$.

Type constructor	T
Class name	C
Type variable	α, β
Type	$\tau ::= \alpha \mid T \mid \tau \tau$
Type constraint	$\pi ::= C \bar{\tau}$
Constraint set	$Q, P ::= \emptyset \mid \pi, Q$
Type scheme	$\sigma ::= \forall \bar{\alpha}. Q \Rightarrow \tau$
Constraint axioms	$\Theta ::= \emptyset \mid \forall \bar{\alpha}. Q \Rightarrow \pi, \Theta$

Fig. 1 Constrained types and context

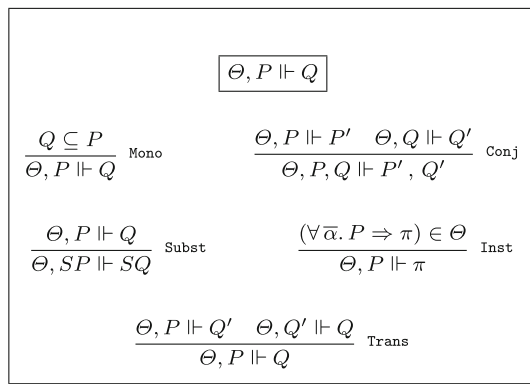


Fig. 2 Type class constraint entailment

Function tv is overloaded, yielding the set of free type variables of types, constraints or constraint sets, and is defined as usual. Sequence $\bar{\alpha}$ used in the context of a set denotes of course the set of type variables in the sequence. The set of constraint axioms Θ is induced by class and instance declarations of a program. Each instance declaration $P \Rightarrow \pi$ where ... introduces an axiom scheme $\forall \bar{\alpha}. P \Rightarrow \pi$, where $\bar{\alpha} = tv(P \Rightarrow \pi)$.

For simplicity and to avoid clutter, in this paper constraint axioms introduced by type class declarations are not considered, since they add no additional problems with respect to termination of constraint set satisfiability and simplification algorithms.

The entailment relation for type class constraints is defined in Fig. 2. Rule (Mono) expresses the property of monotonicity, (Trans) of transitivity, (Subst) of closure under type substitution (cf. [12]), (Inst) defines entailment according to a constraint axiom and (Conj) deals with sets with more than one constraint.

A type substitution S is a (kind-preserving) function from type variables to types, and extends straightforwardly to constraints, and to sets of types and sets of constraints. For convenience, a substitution is often written as a finite mapping $[\alpha_1 \mapsto \tau_1, \dots, \alpha_n \mapsto \tau_n]$, which is also abbreviated as $[\bar{\alpha} \mapsto \bar{\tau}]$. Juxtaposition $S'S$ is used as a synonym for function composition, $S' \circ S$, the domain of a substitution S is defined by $dom(S) = \{\alpha \mid S(\alpha) \neq \alpha\}$ and the restriction of S to V is given by $S|_V(\alpha) = S(\alpha)$ if $\alpha \in V$, otherwise α .

3.1 Constraint set satisfiability

Constraint set satisfiability is central to the interpretation of constrained types and is closely related to simplification and improvement. Following [13], $[P]_\Theta$ denotes the set of satisfiable instances of constraint set P , with respect to constraint axioms Θ :

$$[P]_\Theta = \{SP \mid \Theta \vdash SP\}$$

Equality of constraint sets is considered modulo type variable renaming. That is, constraint sets P and Q are considered to be equal by considering also that a renaming substitution S can be applied to P so as to make SP and Q equal. A substitution S is a renaming substitution if for all $\alpha \in dom(S)$ we have that $S(\alpha) = \beta$, for some type variable $\beta \notin dom(S)$.

If $SP \in [P]_\Theta$ then S is called a satisfying substitution for P .

Subscript Θ will not be used hereafter because satisfiability is always considered with respect to a set of global constraint axioms Θ .

For any substitution S and constraint set P we have that $[SP] \subseteq [P]$. The reverse inclusion, $[P] \subseteq [SP]$, does not always hold, and allow us to characterize improvement of the set of constraints P to an equivalent but simpler or more informative constraint set SP , such that $[SP] = [P]$. Substitution S is called an improving substitution for P if applying S to P preserves the set of satisfiable instances, that is, if $[SP] = [P]$.

The next section presents constraint set satisfiability algorithms, including an algorithm that uses a criterion for guaranteeing termination on any given input. This termination criterion is used in Sect. 5, to define a constraint set simplification algorithm.

4 Computing constraint set satisfiability

Figure 3 presents a computable function that, given any constraint set P , returns, if it terminates, the set of satisfying substitutions for P . The definition uses judgements of the form $\Theta \vdash^{sats} P \rightsquigarrow \mathbb{S}$, meaning that \mathbb{S} is the set of satisfying substitutions for P , with respect to constraint axioms Θ . The following function is used:

$$\begin{aligned} sats(\pi, \Theta) = \{ & (S|_{tv(\pi)}, SP, \pi_0) \mid (\forall \bar{\alpha}. P_0 \Rightarrow \pi_0) \in \Theta, \\ & S_1 = [\bar{\alpha} \mapsto \bar{\beta}], \bar{\beta} \text{ fresh}, \\ & (P \Rightarrow \pi') = S_1 (P_0 \Rightarrow \pi_0), \\ & S = mgu(\pi = \pi') \} \end{aligned}$$

where function mgu gives a most general unifier for a pair of constraints, written as an equality. That is, $mgu(C \bar{\tau} = C \bar{\tau}')$ gives a substitution S such that, $S \bar{\tau} = S \bar{\tau}'$ and, for any S' such that $S' \bar{\tau} = S' \bar{\tau}'$, it holds that $S' = S'' \circ S$, for some S'' .¹

Let \mathbb{S} be the returned set of satisfying substitutions for a given constraint P . Since $S \in \mathbb{S}$ implies $dom(S) \subseteq tv(\pi)$ — because if S is in $sats(\pi, \Theta)$ then $dom(S) \subseteq tv(\pi)$ —the only possible satisfying substitution to be returned for the empty set of constraints is the identity substitution (id), as defined by rule $SEmpty$. Rule $SInst$ computes the set \mathbb{S}_0

¹ See, for example [2], for the general theory of unification and algorithms for computing a most general unifier for a set of term equalities.

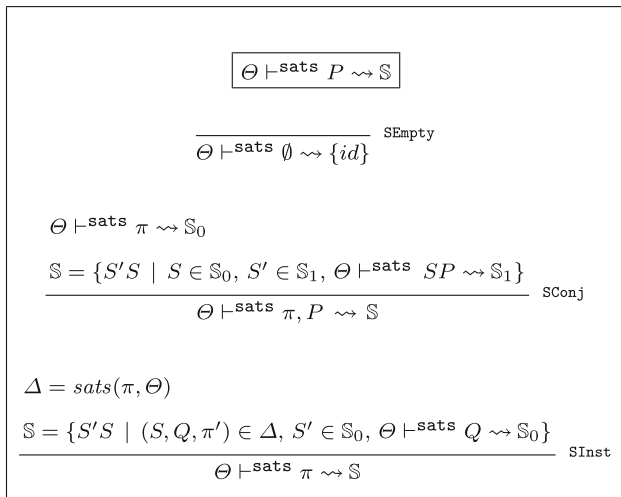


Fig. 3 Constraint set satisfiability

of satisfying substitutions $S \in \mathbb{S}_0$ for a given constraint π , by determining the set of constraint axioms $\forall \bar{\alpha}. P_0 \Rightarrow \pi_0$ in Θ such that π unifies with π_0 , and composing these substitutions with those obtained by recursively computing the set of satisfying substitutions for contexts $S P_0$. Rule $SConj$ deals with sets of constraints. The following examples illustrate the use of these rules.

B , I and F are used in the sequel as abbreviations of *Bool*, *Int* and *Float*, respectively.

Example 1 Consider $P = \{A a b, D b\}$ and

$$\Theta = \{A I [I], A I [B], C I, \forall b. C b \Rightarrow D [b]\}$$

Satisfiability of P with respect to Θ yields a set of substitutions \mathbb{S} given by:

$$\frac{\Theta \vdash^{sats} A a b \rightsquigarrow \mathbb{S}_0 \quad \mathbb{S} = \{S'S \mid S \in \mathbb{S}_0, S' \in \mathbb{S}_1, \Theta \vdash^{sats} S(D b) \rightsquigarrow \mathbb{S}_1\}}{\Theta \vdash^{sats} A a b, \{D b\} \rightsquigarrow \mathbb{S}} SConj$$

Then:

$$\frac{\Delta_0 = \{(S_1, \emptyset, A I [I]), (S_2, \emptyset, A I [B])\} \quad \mathbb{S}_0 = \{S'S \mid (S, Q, \pi') \in \Delta_0, S' \in \mathbb{S}', \Theta \vdash^{sats} Q \rightsquigarrow \mathbb{S}'\}}{\Theta \vdash^{sats} A a b \rightsquigarrow \mathbb{S}_0} SInst$$

where $S_1 = [a \mapsto I, b \mapsto [I]]$, $S_2 = [a \mapsto I, b \mapsto [B]]$.

Then, by rule $SConj$, the set of satisfying substitutions for $S_1(D b) = D [I]$ and $S_2(D b) = D [B]$ must be computed, and are given respectively by:

$$\frac{\Delta_1 = \{(S'_1 | \emptyset, \{C I\}, D [b])\} \quad \mathbb{S}_1^1 = \{S'S \mid (S, Q, \pi') \in \Delta_1, S' \in \mathbb{S}', \Theta \vdash^{sats} Q \rightsquigarrow \mathbb{S}'\}}{\Theta \vdash^{sats} D [I] \rightsquigarrow \mathbb{S}_1^1} SInst$$

where $S'_1 = [b_1 \mapsto I]$, b_1 is a fresh type variable, $S'_1 | \emptyset = id$, and

$$\Delta_2 = \{(S'_2 | \emptyset, \{C B\}, D [b])\} \quad \mathbb{S}_1^2 = \{S'S \mid (S, Q, \pi') \in \Delta_2, S' \in \mathbb{S}', \Theta \vdash^{sats} Q \rightsquigarrow \mathbb{S}'\}$$

$$\frac{}{\Theta \vdash^{sats} D [B] \rightsquigarrow \mathbb{S}_1^2} SInst$$

where $S'_2 = [b_2 \mapsto B]$, b_2 is a fresh type variable, $S'_2 | \emptyset = id$. Now, $\mathbb{S}_1^1 = \{id\}$ and $\mathbb{S}_1^2 = \emptyset$. Thus, $\mathbb{S} = \{S_1\}$.

The example below, extracted from [3], illustrates non-termination of the computation of the set of satisfying substitutions by the function defined in Fig. 3. We use $T^2 \tau$ to abbreviate $T(T \tau)$ and similarly for other indices greater than 2.

Example 2 Let $\Theta = \{\forall a, b. \{C a b\} \Rightarrow C (T^2 a) b\}$ and consider computing satisfiability of $\pi = C a (T a)$ with respect to Θ .

We have that π unifies with the head of constraint axiom $\forall a, b. (C a b) \Rightarrow C (T^2 a) b$, giving substitution $S = [a \mapsto T^2 a_1, b_1 \mapsto T^3 a_1]$. We must then recursively compute the set of satisfying substitutions of constraint $S(C a_1 b_1) = C a_1 (T^3 a_1)$. This constraint also unifies with $\forall a, b. (C a b) \Rightarrow C (T^2 a) b$, giving substitution $S_1 = [a_1 \mapsto (T^2 a_2), b_2 \mapsto (T^3 a_1 = T^5 a_2)]$. Again, we must recursively compute the set of satisfying substitutions of constraint $S_1(C a_2 b_2) = C a_2 (T^5 a_2)$, and the process goes on forever.

The following theorems state, respectively, correctness and completeness of the constraint set satisfiability algorithms presented in Fig. 3, with respect to the entailment relation.

Theorem 1 (Correctness of \vdash^{sats}) If $\Theta \vdash^{sats} P \rightsquigarrow \mathbb{S}$ then $\Theta \Vdash S P$, for all $S \in \mathbb{S}$.

Proof By induction over the derivation of $\Theta \vdash^{sats} P \rightsquigarrow \mathbb{S}$. The only interesting case is for rule $SInst$. Let $\pi = C \bar{\tau}$ and $\Delta = sats(\pi, \Theta)$. If $\Delta = \emptyset$, the theorem holds trivially. Thus, assume $\Delta \neq \emptyset$ and let $(S, Q, C \bar{\tau}_0) \in \Delta$. By the definition of $sats$, this means that $\forall \bar{\alpha}. P_0 \Rightarrow C \bar{\tau}_0 \in \Theta$, where $\bar{\alpha} = tv(P_0 \Rightarrow C \bar{\tau}_0)$, and $P' \Rightarrow C \bar{\tau} = [\bar{\alpha} \mapsto \bar{\beta}] P_0 \Rightarrow C \bar{\tau}_0$. By rule $Inst$ we have that $\Theta, P_0 \Vdash C \bar{\tau}_0$ is provable. We also have that $\Theta \vdash^{sats} Q \rightsquigarrow \mathbb{S}_0$, where $Q = S[\bar{\alpha} \mapsto \bar{\beta}] P_0$, and thus, by the induction hypothesis, we have that (1) $\Theta \Vdash S' Q$ holds for all $S' \in \mathbb{S}_0$. Also, since $\Theta, P_0 \Vdash C \bar{\tau}_0$ is provable, we have, by rule $Subst$, that (2) $\Theta, S_0 P_0 \Vdash S_0 C \bar{\tau}_0$, where $S_0 = S' S[\bar{\alpha} \mapsto \bar{\beta}]$. From (1) and (2) we have, by rule $Trans$, that $\Theta \Vdash S_0 C \bar{\tau}_0$ is provable. Since $S \bar{\tau} = S[\bar{\alpha} \mapsto \bar{\beta}] \bar{\tau}_0$, this means that $\Theta \Vdash S' S C \bar{\tau}$ is provable. \square

Theorem 2 (Completeness of \vdash^{sats}) If $\Theta \Vdash S P$ then there exist $S' \in \mathbb{S}$ and S'' such that $S'' S' P = S P$, where $\Theta \vdash^{sats} P \rightsquigarrow \mathbb{S}$.

Proof Induction over $S P$ in $\Theta \Vdash S P$.

4.1 Termination

The algorithm presented in Fig. 3 is modified in this section in order to ensure termination on any given input. The basic idea is to associate a value to each constraint head of the set of constraint axioms that is unified with some constraint in the recursive process of computing satisfiability, and require that the value associated to a constraint head always decreases in a new unification that occurs during this process. Computation stops if this requirement is not fulfilled, with no satisfying substitution found for the original set of constraints. Values in this decreasing chain are a measure of the size of types in constraints that unify with each constraint head axiom: the size of each constraint in this chain is decreasing or there exists a position of a type argument in the constraint such that the type's size is decreasing.

Let the constraint value $\eta(\pi)$ of a constraint π , which gives the number of occurrences of type variables and type constructors in π , be defined as follows:

$$\begin{aligned}\eta(C \tau_1 \cdots \tau_n) &= \sum_{i=1}^n \eta(\tau_i) \\ \eta(T) &= 1 \\ \eta(\alpha) &= 1 \\ \eta(\tau_1 \tau_2) &= \eta(\tau_1) + \eta(\tau_2)\end{aligned}$$

A finite constraint-head-value function Φ is used to map constraint heads π_0 of Θ to pairs (I, Π) , as follows.

The first component I is a tuple (v_0, \dots, v_n) , where v_0 is the least $\eta(S\pi')$ of all constraints π' that have unified with π_0 during the satisfiability test for π , where $S = \text{mgu}(\pi'_0, \pi')$. Each v_i , $1 \leq i \leq n$, is the least $\eta(\tau_i)$ where τ_i is a type belonging to some $S\pi'$ that has unified with π_0 .

We let $I.v_i$ denote the i -th value of I and, similarly, $\Phi(\pi_0).I$ and $\Phi(\pi_0).\Pi$ denote respectively the first and second components of $\Phi(\pi_0)$.

The second component Π of $\Phi(\pi_0)$ contains constraints π' that unify with π_0 and have constraint values equal to v_0 . This allows distinct constraints with equal constraint values to unify with π_0 (cf. Example 6 below).

Consider a recursive step in a test of satisfiability where a constraint π unifies with a constraint head $\pi_0 = C \tau_1 \dots \tau_n$, with $S = \text{mgu}(\pi_0, \pi)$. Let $\Phi(\pi_0) = ((v_0, \dots, v_n), \Pi)$ and $\eta(S\pi) = n_0$. $\Phi(\pi_0)$ is then updated as follows. If $n_0 < v_0$ then only the value v_0 is updated, to n_0 . In the case that $n_0 = v_0$ and $\pi \notin \Pi$, $\Phi(\pi_0)$ is updated to $((v_0, \dots, v_n), \Pi \cup \{S\pi\})$, i.e. we include $S\pi$ in the set of constraints that have the same value v_0 . Finally, if $n_0 > v_0$, we set v_0 to -1 and for each τ_i such that $\eta(\tau_i) \geq v_i$, we update v_i with -1 , otherwise v_i is updated with $\eta(\tau_i)$. In subsequent steps for constraints π' that unify with π_0 , with S' as a unifying substitution, it is required that $\eta(S'\tau_i) < v_i$; if there's no such i , a failure in the termination criteria is detected.

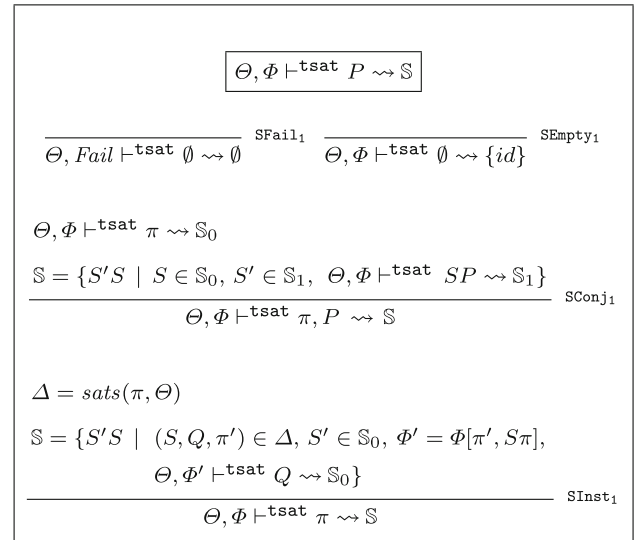


Fig. 4 Terminating constraint set satisfiability

Let $f[x \mapsto y]$ denote the usual function updating notation for f' given by $f'(x') = y$ if $x' = x$, otherwise $f(x)$.

We define $\Phi[\pi_0, \pi]$ as updating of $\Phi(\pi_0) = (I, \Pi)$ as follows, where $I = (v_0, v_1, \dots, v_n)$, $\pi = C \tau_1 \cdots \tau_n$, $n_0 = \eta(\pi)$:

$$\begin{aligned}\Phi[\pi_0, \pi] &= \Phi[\pi_0 \mapsto ((n_0, v_1, \dots, v_n), \Pi)] \text{ if } n_0 < I.v_0; \\ &\Phi[\pi_0 \mapsto (I, \Pi \cup \{\pi\})] \text{ if } n_0 = I.v_0, \pi \notin \Pi; \\ &\Phi[\pi_0 \mapsto (I', \Pi)] \text{ if } n_0 > I.v_0, \exists i. (I'.v_i \neq -1) \\ &\text{where, for } i = 1, \dots, n, \\ &I'.v_i = \begin{cases} -1 & \text{if } I.v_i < \eta(\tau_i) \text{ or } i = 0 \\ \eta(\tau_i) & \text{otherwise} \end{cases} \\ &\text{Fail otherwise} \end{aligned}$$

The computable function (tsat) for constraint satisfiability, defined in Fig. 4, uses judgements of the form $\Theta, \Phi \vdash^{\text{tsat}} P \rightsquigarrow \mathbb{S}$, with constraint-head-value function Φ as additional parameter.

The set of satisfying substitutions for constraint set P with respect to the set of constraint axioms Θ is given by \mathbb{S} , such that $\Theta, \Phi_0 \vdash^{\text{tsat}} P \rightsquigarrow \mathbb{S}$ holds, where $\Phi_0(\pi_0) = (I_0, \emptyset)$ for each constraint head $\pi_0 = C \tau_1 \dots \tau_n$ in Θ and I_0 is a tuple formed by $n + 1$ occurrences of a large enough integer constant, represented by ∞ .

Consider the following.

Example 3 Consider computing satisfiability of $\pi = \text{Eq}[[I]]$ in $\Theta = \{\text{Eq}I, \forall a. \text{Eq}a \Rightarrow \text{Eq}[a]\}$, letting $\pi_0 = \text{Eq}[a]$; we have:

$$\begin{aligned}\Delta_0 &= \text{sats}(\pi, \Theta) = \{(S|_{\emptyset}, \{\text{Eq}[I]\}, \pi_0)\} \\ S &= [a_1 \mapsto [I]] \\ \mathbb{S}_0 &= \{S_1 \circ id \mid S_1 \in \mathbb{S}_1, \Theta, \Phi_1 \vdash^{\text{tsat}} \text{Eq}[I] \rightsquigarrow \mathbb{S}_1\} \\ &\Theta, \Phi_0 \vdash^{\text{tsat}} \pi \rightsquigarrow \mathbb{S}_0\end{aligned}$$

where $\Phi_1 = \Phi_0[\pi_0, \pi]$, $\Phi_1(\pi_0).I = (\eta(\pi) = 3, \infty)$, $S\pi = \pi$ and a_1 is a fresh type variable; then:

$$\begin{array}{l} \Delta_1 = \text{sats}(Eq[I], \Theta) = \{(S' | \emptyset, \{Eq[I], \pi_0\})\} \\ S' = [a_2 \mapsto I] \\ \mathbb{S}_1 = \{S_2 \circ id \mid S_2 \in \mathbb{S}_2, \Theta, \Phi_2 \vdash^{\text{tsat}} Eq[I] \rightsquigarrow S_2\} \\ \hline \Theta, \Phi_1 \vdash^{\text{tsat}} Eq[I] \rightsquigarrow \mathbb{S}_1 \end{array}$$

where $\Phi_2 = \Phi_1[\pi_0, Eq[I]]$ and $\eta(Eq[I]) = 2$ is less than $\Phi_1(\pi_0).I.v_0 = 3$; then:

$$\begin{array}{l} \Delta_2 = \text{sats}(Eq[I], \Theta) = \{(id, \emptyset, Eq[I])\} \\ \mathbb{S}_2 = \{S_3 \circ id \mid S_3 \in \mathbb{S}_3, \Theta, \Phi_3 \vdash^{\text{tsat}} \emptyset \rightsquigarrow S_3 = \{id\}\} \\ \hline \Theta, \Phi_2 \vdash^{\text{tsat}} Eq[I] \rightsquigarrow \mathbb{S}_2 \end{array}$$

where $\Phi_3 = \Phi_2[Eq[I], Eq[I]]$, $\mathbb{S}_3 = \{id\}$ by (SEmpty_Y).

Example 4 Consider again Example 2: we want to obtain the set of satisfying substitutions for constraint $\pi = C a (T a)$, given $\Theta = \{\forall a, b. C a b \Rightarrow C (T^2 a) b\}$ (computation with input π by the function in Fig. 3 does not terminate). We have, where $\pi_0 = C (T^2 a) b$:

$$\begin{array}{l} \Delta_0 = \text{sats}(\pi, \Theta) = \{(S' |_{\{a\}}, \{\pi_1\}, \pi_0)\} \\ S' = [a \mapsto T^2 a_1, b_1 \mapsto T^3 a_1] \\ \pi_1 = C a_1 (T^3 a_1) \\ \mathbb{S}_0 = \{S_1 \circ [a \mapsto T^2 a_1] \mid S_1 \in \mathbb{S}_1, \Theta, \Phi_1 \vdash^{\text{tsat}} \pi_1 \rightsquigarrow \mathbb{S}_1\} \\ \hline \Theta, \Phi_0 \vdash^{\text{tsat}} \pi \rightsquigarrow \mathbb{S}_0 \end{array}$$

where $\Phi_1 = \Phi_0[\pi_0, S\pi]$, $\eta(S\pi) = \eta(C (T^2 a_1) (T^3 a_1)) = 7 < \Phi_0(\pi_0).I.v_0 = \infty$; then:

$$\begin{array}{l} \Delta_1 = \text{sats}(\pi_1, \Theta) = \{(S' |_{\{a_1\}}, \{\pi_2\}, \pi_0)\} \\ S' = [a_1 \mapsto T^2 a_2, b_2 \mapsto T^3 a_1 = T^5 a_2] \\ \pi_2 = C a_2 (T^5 a_2) \\ \mathbb{S}_1 = \{S_2 \circ [a_1 \mapsto T^2 a_2] \mid S_2 \in \mathbb{S}_2, \Theta, \Phi_2 \vdash^{\text{tsat}} \pi_2 \rightsquigarrow \mathbb{S}_2\} \\ \hline \Theta, \Phi_1 \vdash^{\text{tsat}} \pi_1 \rightsquigarrow \mathbb{S}_1 \end{array}$$

where $\Phi_2 = \Phi_1[\pi_0, S'\pi_1]$, $S'\pi_1 = (C (T^2 a_2) (T^5 a_2))$ and, since $\eta(S'\pi_1) = 9 > \Phi_1(\pi_0).I.v_0 = 7$, we have that $\Phi_2(\pi_0).I = (-1, \eta(T^2 a_2) = 3, \eta(T^5 a_2) = 6)$; then:

$$\begin{array}{l} \Delta_1 = \text{sats}(\pi_2, \Theta) = \{(S' |_{\{a_2\}}, \{\pi_3\}, \pi_0)\} \\ S'' = [a_2 \mapsto T^2 a_3, b_3 \mapsto T^5 a_2 = T^7 a_3] \\ \pi_3 = C a_3 (T^7 a_3) \\ \mathbb{S}_2 = \{S_3 \circ [a_2 \mapsto T^2 a_3] \mid S_3 \in \mathbb{S}_3, \Theta, \Phi_3 \vdash^{\text{tsat}} \pi_3 \rightsquigarrow \mathbb{S}_3\} \\ \hline \Theta, \Phi_1 \vdash^{\text{tsat}} \pi_1 \rightsquigarrow \mathbb{S}_2 \end{array}$$

where $\Phi_3 = \Phi_2[\pi_0, S''\pi_2] = \text{Fail}$, because $\eta(S''\pi_2) = \eta(C (T^3 a_3) (T^7 a_3)) = 12 > \Phi_2(\pi_0).I.v_0 = 9$ and there's no i such that $\Phi_3(\pi_0).I.v_i \neq -1$, meaning that no parameter of $S''\pi_2$ has a decreasing η value.

The following illustrates an example of a satisfiable constraint for which computation of satisfiability involves computing satisfiability of constraints π' that unify with a constraint head π_0 such that $\eta(\pi')$ is greater than the upper bound associated to π_0 .

Example 5 Consider satisfiability of $\pi = C \mathbb{I} (T^3 \mathbb{I})$ in $\Theta = \{C (T a) \mathbb{I}, \forall a, b. C (T^2 a) b \Rightarrow C a (T b)\}$. We have, where $\pi_0 = C a (T b)$:

$$\begin{array}{l} \Delta_0 = \text{sats}(\pi, \Theta) = \{(S' | \emptyset, \{\pi_1\}, \pi_0)\} \\ S' = [a_1 \mapsto \mathbb{I}, b_1 \mapsto T^2 \mathbb{I}] \\ \pi_1 = C (T^2 \mathbb{I}) (T^2 \mathbb{I}) \\ \mathbb{S}_0 = \{S_1 \circ id \mid S_1 \in \mathbb{S}_1, \Theta, \Phi_1 \vdash^{\text{tsat}} \pi_1 \rightsquigarrow \mathbb{S}_1\} \\ \hline \Theta, \Phi_0 \vdash^{\text{tsat}} \pi \rightsquigarrow \mathbb{S}_0 \end{array}$$

where $\Phi_1 = \Phi_0[\pi_0, \pi]$, $\eta(\pi) = 5 < \Phi_0(\pi_0).I.v_0 = \infty$, $S\pi = \pi$; then:

$$\begin{array}{l} \Delta_1 = \text{sats}(\pi_1, \Theta) = \{(S' | \emptyset, \{\pi_2\}, \pi_0)\} \\ S' = [a_2 \mapsto T^2 \mathbb{I}, b_2 \mapsto T \mathbb{I}] \\ \pi_2 = C (T^4 \mathbb{I}) (T \mathbb{I}) \\ \mathbb{S}_1 = \{S_2 \circ [a_1 \mapsto T^2 a_2] \mid S_2 \in \mathbb{S}_2, \Theta, \Phi_2 \vdash^{\text{tsat}} \pi_2 \rightsquigarrow \mathbb{S}_2\} \\ \hline \Theta, \Phi_1 \vdash^{\text{tsat}} \pi_1 \rightsquigarrow \mathbb{S}_1 \end{array}$$

where $\Phi_2 = \Phi_1[\pi_0, \pi_1]$, $\Phi_1(\pi_0).I = (5, \infty, \infty)$, $S'\pi_1 = \pi_1$. Since $\eta(\pi_1) = 6 > 5 = \Phi_1(\pi_0).I.v_0$, we have that $\Phi_2(\pi_0).I$ becomes equal to $(-1, 3, 3)$.

Then, consider that $\pi_2 = C \tau_1 \tau_2$ where $\tau_1 = T^4 I$ and $\tau_2 = T I$. Since $\eta(\pi_2) > \Phi_2(\pi_0).I.v_0 = -1$, there must exist i , $1 \leq i \leq 2$, such that $\eta(\tau_i) < \Phi_2(\pi_0).v_i$, and such condition is satisfied for $i = 2$, updating $\Phi_2(\pi_0).I$ to $(-1, -1, 2)$. Satisfiability is then finally tested for $\pi_3 = C (T^6 \mathbb{I}) \mathbb{I}$, that unifies with $\pi_0 = C (T a) \mathbb{I}$, which returns $\mathbb{S}_3 = \{[a_3 \mapsto T^5 \mathbb{I}] | \emptyset\} = \{id\}$. Constraint π is thus satisfiable, with $\mathbb{S}_0 = \{id\}$.

The following example illustrates the use of a set of constraints as a component of the constraint-head-value function.

Example 6 Let $\pi = C (T^2 \mathbb{I}) F$, $\pi_0 = C (T a) b$, $\Theta = \{C \mathbb{I} (T^2 F), \forall a, b. C a (T b) \Rightarrow C (T a) b\}$:

$$\begin{array}{l} \Delta_0 = \text{sats}(\pi, \Theta) = \{(S' | \emptyset, \{\pi_1\}, \pi_0)\} \\ S' = [a_1 \mapsto (T \mathbb{I}), b_1 \mapsto F], \pi_1 = C (T \mathbb{I}) (T F) \\ \mathbb{S}_0 = \{S_1 \circ id \mid S_1 \in \mathbb{S}_1, \Theta, \Phi_1 \vdash^{\text{tsat}} \pi_1 \rightsquigarrow \mathbb{S}_1\} \\ \hline \Theta, \Phi_0 \vdash^{\text{tsat}} \pi \rightsquigarrow \mathbb{S}_0 \end{array}$$

where $\Phi_1 = \Phi_0[\pi_0, \pi]$, $S\pi = \pi$; then:

$$\begin{array}{l} \Delta_1 = \text{sats}(\pi_1, \Theta) = \{(S' | \emptyset, \{\pi_2\}, \pi_0)\} \\ S' = [a_2 \mapsto \mathbb{I}, b_2 \mapsto T F], \pi_2 = C \mathbb{I} (T^2 F) \\ \mathbb{S}_1 = \{S_2 \circ id \mid S_2 \in \mathbb{S}_2, \Theta, \Phi_2 \vdash^{\text{tsat}} \pi_2 \rightsquigarrow \mathbb{S}_2\} \\ \hline \Theta, \Phi_1 \vdash^{\text{tsat}} \pi_1 \rightsquigarrow \mathbb{S}_1 \end{array}$$

where $\Phi_2 = \Phi_1[\pi_0, \pi_1]$, $\eta(\pi_1) = 4 = \Phi_1(\pi_0).I.v_0 = \eta(\pi)$, $S'\pi_1 = \pi_1$ and π_1 is not in $\Phi_1(\pi_0). \Pi_1 = \emptyset$. We have that $\mathbb{S}_2 = \{id\}$, because $\text{sats}(C \mathbb{I} (T^2 F), \Theta) = \{(id, \emptyset, C \mathbb{I} (T^2 F))\}$, and π is then satisfiable.

Since satisfiability of type class constraints is in general undecidable [6], there exist instances of this problem for which our algorithm incorrectly reports unsatisfiability. An example that exhibits an incorrect behavior, constructed by encoding a solvable post correspondence problem

(PCP) instance by means of constraint set satisfiability, using G. Smith's scheme [6], is shown below. For all examples mentioned in the literature [15, 17] and numerous tests that include those used by GHC involving pertinent GHC extensions, the algorithm works as expected, without the need of any compilation flag.

Example 7 This example uses a PCP instance taken from [9]. A PCP instance can be defined as composed of pairs of strings, each pair having a top and a bottom string, where the goal is to select a sequence of pairs such that the two strings obtained by concatenating top and bottom strings in such pairs are identical. The example uses three pairs of strings: $p_1 = (100, 1)$ (that is, pair 1 has string 100 as the top string and 1 as the bottom string), $p_2 = (0, 100)$ and $p_3 = (1, 00)$.

This instance has a solution: using numbers to represent corresponding pairs (i.e., 1 represents pair 1 and analogously for 2 and 3), the sequence of pairs 1311322 is a solution.

A satisfiability problem that has a solution if and only if the PCP instance has a solution can be constructed by adapting G. Smith's scheme to Haskell's notation. We consider for this a two-parameter class C , and a constraint context such that $\Theta = \Theta_1 \cup \Theta_2 \cup \Theta_3$, where Θ_i is constructed from pair i , for $i = 1, 2, 3$:

$$\begin{aligned}\Theta_1 &= \{ C(1 \rightarrow 0 \rightarrow 0)1, \\ &\quad \forall a, b. C a b \Rightarrow C(1 \rightarrow 0 \rightarrow 0 \rightarrow a)(1 \rightarrow b) \} \\ \Theta_2 &= \{ C 0(1 \rightarrow 0 \rightarrow 0), \\ &\quad \forall a, b. C a b \Rightarrow C(0 \rightarrow a)(1 \rightarrow 0 \rightarrow 0 \rightarrow b) \} \\ \Theta_3 &= \{ C 1(0 \rightarrow 0), \\ &\quad \forall a, b. C a b \Rightarrow C(1 \rightarrow a)(0 \rightarrow 0 \rightarrow b) \}\end{aligned}$$

We have that constraint $C a a$ is satisfiable, with a solution constructed from solution 1311322 of the PCP instance. Computation by our algorithm terminates, erroneously reporting unsatisfiability. The steps of the computation are omitted. The error occurs because a constraint $\pi_2 = C a_2(1 \rightarrow a_2)$ unifies with $\pi_{01} = C(1 \rightarrow 0 \rightarrow 0 \rightarrow a)(1 \rightarrow b)$ and $\eta(S\pi_2)$ is greater than $\Phi(\pi_{01}).I.v_0$, where $S = mgu(\pi_2, \pi_{01})$, and there's no $i \in \{1, 2\}$ such that $\Phi_3(\pi_0).I.v_i \neq -1$, meaning that no parameter of $S\pi_2$ has a decreasing η value.

To prove that the computation of the set of satisfying substitutions for any given constraint set P by the function defined in Fig. 4 always terminates, consider that an infinite recursion might only occur if an infinite number of constraints unified with the head π_0 of one constraint axiom in Θ , since there exist finitely many constraint axioms in Θ . This is avoided because, for any new constraint π that unifies with π_0 , we have, by the definition of $\Phi[\pi_0, \pi]$, that $\Phi(\pi_0)$ is updated to a value distinct from the previous ones (otherwise $\Phi[\pi_0, \pi]$ yields *Fail* and computation is stopped). The conclusion follows from the fact that $\Phi(\pi_0)$ can have only

finitely many distinct values, for any π_0 . This can be seen by considering that, for any π_0 such that $\Phi(\pi_0) = (I, \Pi)$, the insertion of a new constraint in Π decreases $k - k'$, where k is the finite number of all possible values that can be inserted in Π and k' is the cardinality of Π . Such a decrease causes then a decrease of Φ (since there exists only finitely many constraint heads π_0 in Θ). Similarly, at each step there must exist some i such that $I.v_i$ decreases, and this can happen only a finitely number of times. We conclude that computation on any given input terminates.

The proposed termination criteria is related to the *Paterson Condition* used in the GHC compiler (see Sect. 2). The constraint value is based on item 2 of this condition, but, instead of using it as a syntactic restriction over constraint heads and contexts in instance declarations, we use it in the definition of a finitely decreasing chain over recursively dependent constraints.

In comparison to the use of a recursion depth limit, our approach has the advantage that type-correctness is not implementation dependent (a constraint is or is not satisfiable with respect to a given set of constraint axioms). The use of a recursion depth limit can make a constraint set satisfiable in one implementation and unsatisfiable in another that uses a lower limit. Incorrectly reporting unsatisfiability can occur in both cases, but is expected to be extremely rare with our approach. We are not aware of any practical example where this occurs.

The main disadvantages of our approach are that it is not syntactically possible to characterize such incorrect unsatisfiability cases and it is not very easy for programmers to understand how type class constraints are handled in such a case, if and when it occurs. However, we expect these cases not to occur in practice.

The presented algorithm has been verified to behave correctly, without the need of any compilation flag, on all examples found in the literature [15], all GHC test cases, involving flags `FlexibleInstances`, `FlexibleContexts` and `UndecidableInstances`, and on Haskell libraries that use multi-parameter type classes, including the monad transformer library [1].

5 Constraint set simplification

The process of simplification of a constraint set, also called context reduction, consists of reducing each constraint π in this set to the context obtained by recursively reducing the context P of the *matching instance* for π in Θ , if such matching exists, until $P = \emptyset$ or there exists no instance in Θ that matches with π . In the latter case π reduces to itself.

This recursive process may not terminate: as a simple example, consider reduction of constraint $C a$ when $\Theta = \{\forall a. C a \Rightarrow C a\}$.

This section presents a computable function for constraint set simplification, where computation is guaranteed to terminate by using the same criterion used in Sect. 4.1.

Constraint set simplification is essentially based on instance matching. We use function $matches(\pi, \Theta)$, defined below, in order to capture the relevant information of matching constraint axioms in Θ with a given constraint π . Function $matches$ is defined by using function $sats$ (Sect. 4), through skolemization of type variables that occur in the given constraint argument (Skolem variables are non unifiable variables, that is, constants):

$$matches(\pi, \Theta) = \left\{ (S P, \pi') \mid \Delta = sats([\bar{\alpha} \mapsto \bar{K}]\pi, \Theta), \right. \\ \left. (S, S P, \pi') \in \Delta, \bar{\alpha} = tv(\pi), \right. \\ \left. \bar{K} \text{ are fresh Skolem variables} \right\}$$

Function $matches(\pi, \Theta)$ returns either a singleton or an empty set².

Constraint set simplification uses a function defined in Fig. 5 by means of judgements of the form $\Theta, \Phi \vdash^{simp} P \rightsquigarrow Q$. This means that reduction of constraint set P under constraint axioms Θ either give constraint set Q as a result or fails. Failure is caused by the criterion used for ensuring termination, explained in Sect. 4.1. Using this function, context reduction is defined as follows, where Φ_0 is as defined in Sect. 4.1:

$$\frac{\text{for } i = 1, \dots, n, Q_i = \begin{cases} \pi_i & \text{if } \Theta, \Phi_0 \vdash^{simp} \pi_i \rightsquigarrow fail \\ Q'_i & \text{if } \Theta, \Phi_0 \vdash^{simp} \pi_i \rightsquigarrow Q'_i \end{cases}}{\Theta \vdash^{simp_0} \{\pi_1, \dots, \pi_n\} \rightsquigarrow Q_1, \dots, Q_n} R_0$$

The rules of Fig. 5 are analogous to the ones in Fig. 4, but now termination enforced by the termination criterion is reported as a failure, which must be propagated backwards along the recursive calls of the computation. Thus, reduction of a constraint π is now defined by two rules, (RInst₁) and (RInst₂) and, analogously, two different rules are used for specifying reduction of a non-singleton set of constraints.

Rule (REmpty) specifies that an empty set of constraints reduces to itself. Rule (RStop) specifies that a constraint π cannot be reduced if there is no instance in Θ that matches with π . Rule (RFail) enforces termination, expressing that reduction cannot be performed since updating of Φ fails.

The process of constraint set simplification is illustrated by the following example.

² We do not consider *overlapping instances* [20], since the subject is unrelated to termination of constraint set satisfiability and simplification. Supporting overlapping instances would need a modification of function $matches$ so as to select a single instance if there exist overlapping matching instances.

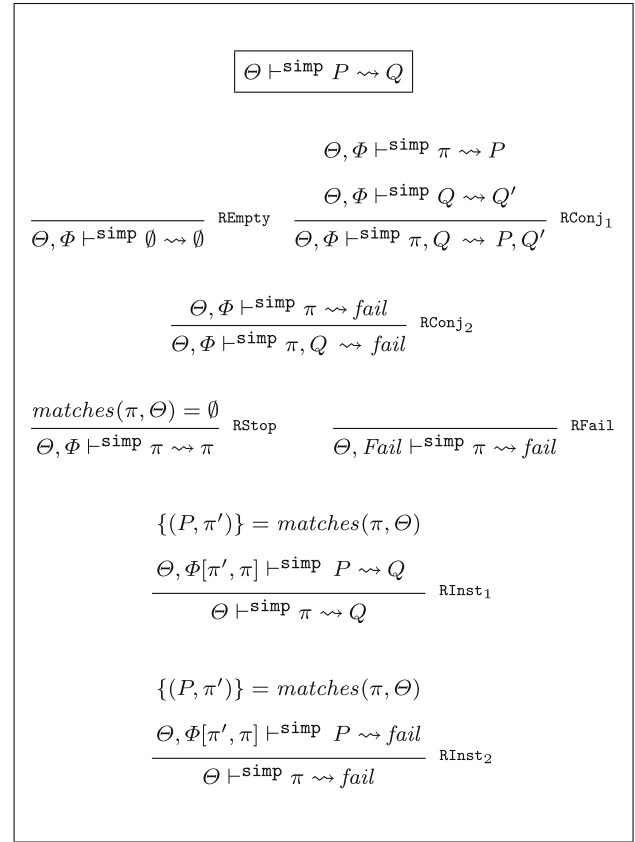


Fig. 5 Constraint set simplification

Example 8 Let $\Theta = \{\forall a. C(T a) \Rightarrow C a, D \top\}$ and $P = \{D \top, C a\}$. According to rule (R₀) reduction of P amounts to independently reducing constraints $D \top$ and $C a$.

Reduction of $D \top$ is defined by rule (RInst₁):

$$\frac{\{(\emptyset, D \top)\} = matches(D \top, \Theta) \quad \Theta, \Phi_0[D \top, D \top] \vdash^{simp} \emptyset \rightsquigarrow \emptyset}{\Theta, \Phi_0 \vdash^{simp} D \top \rightsquigarrow \emptyset}$$

Reduction of $\pi = \pi_0 = C a$ results in failure, as shown below:

$$\frac{\{(C(T a_1), \pi_0)\} = matches(\pi, \Theta) \quad \Theta, \Phi_1 \vdash^{simp} (C(T a_1)) \rightsquigarrow fail}{\Theta, \Phi_0 \vdash^{simp} \pi \rightsquigarrow fail}$$

where $\Phi_1 = \Phi_0[\pi, \pi_0]$, $\Phi_1(\pi_0).I = (\eta(\pi) = 1, \infty)$. We have that:

$$\frac{\{(C(T^2 a_2), \pi_0)\} = matches(C(T a_1), \Theta) \quad \Theta, \Phi_2 \vdash^{simp} (C(T^2 a_2)) \rightsquigarrow fail}{\Theta, \Phi_1 \vdash^{simp} (C(T a_1)) \rightsquigarrow fail}$$

where $\Phi_2 = \Phi_1[C(T a_1), \pi_0] = fail$ because $\eta(C(T a_1)) \not\prec \Phi_1(\pi_0).I.v_1 = 1$.

By rule (R₀), we have that $\Theta \vdash^{simp_0} \{D \top, C a\} \rightsquigarrow \{C a\}$, meaning that $D \top$ can be removed and $C a$ cannot be further reduced.

The following theorem states the correctness of the constraint simplification function defined in Fig. 5.

Theorem 3 [Correctness of \vdash^{simp}] *If $\Theta, \Phi \vdash^{\text{simp}} P \rightsquigarrow Q$ holds, then $\Theta, Q \Vdash P$ is provable and Q cannot be further simplified, i.e., $\Theta, \Phi \vdash^{\text{simp}} Q \rightsquigarrow Q$.*

Proof Induction over $\Theta, \Phi \vdash^{\text{simp}} P \rightsquigarrow Q$.

6 Conclusion

This paper presents a termination criterion and terminating algorithms for constraint simplification and improvement, based on the use of a value that always decreases on each recursive step in these algorithms. The termination criterion defined can be used in any form of constraint simplification and improvement algorithm during type inference.

The use of this criterion eliminates the need for imposing syntactic conditions on Haskell type class and instance declarations and the need for using a recursion stack depth limit in order to guarantee termination of type inference in the presence of multi-parameter type classes, in case these syntactic conditions are chosen by programmers not to be enforced.

Since type class constraint satisfiability is in general undecidable, there exist instances of this problem for which the algorithm presented in this paper incorrectly reports unsatisfiability. However, practical examples where this occurs are expected to be very rare. The algorithms have been implemented and used in a prototype front-end for Haskell, available at <http://github.com/rodrigogribeiro/mptc>. For all examples mentioned in the literature, Haskell libraries that use multi-parameter type classes and tests used by the Haskell GHC compiler, involving all pertinent GHC extensions, the algorithm works as expected without the need for any compilation flag.

In comparison to the use of a recursion depth limit, our approach has the advantage that type-correctness is not implementation dependent (a constraint is or is not satisfiable with respect to a given set of constraint axioms). The use of a recursion depth limit can make a constraint set satisfiable in one implementation and unsatisfiable in another that uses a lower limit. Incorrectly reporting unsatisfiability can occur in both cases, but is expected to be extremely rare with our approach. We are not aware of any practical example where this occurs.

The main disadvantages of our approach are that it is not syntactically possible to characterize such incorrect unsatisfiability cases and it is not very easy for programmers to understand how type class constraints are handled in such a case, if and when it occurs.

Acknowledgments We would like to thank the anonymous reviewers for their careful work, which has been very useful to improve the paper.

References

- Gill A (2006) MTL—The Monad Transformer Library. <http://hackage.haskell.org/package/mtl>
- Baader F, Snyder W (2001) Unification theory. In: Robinson J., Voronkov A (eds) Handbook of Automated Reasoning, Elsevier Science Publishers, vol. 1, pp 447–533
- Camarão C, Figueiredo L, Vasconcellos C (2004) Constraint-set Satisfiability for Overloading. In: Proc. of the 6th ACM SIGPLAN International Conf. on Principles and Practice of Declarative Programming (PPDP'04), pp 67–77
- Camarão C, Ribeiro R, Figueiredo L, Vasconcellos C (2009) A Solution to Haskell's Multi-Parameter Type Class Dilemma. In: Proc. of the 13th Brazilian Symposium on Programming Languages (SBLP'2009), pp 5–18. <http://www.dcc.ufmg.br/camarao/CT/solution-to-mptc-dilemma.pdf>
- Hall C, Hammond K, Jones SP, Wadler P (1996) Type Classes in Haskell. ACM Trans Program Lang Syst 18(2):109–138
- Smith G (1991) Polymorphic type inference for languages with overloading and subtyping. Ph.D. thesis, Cornell Univ.
- Jones M, Diatchki I (2008) Language and Program Design for Functional Dependencies. In: ACM SIGPLAN Haskell, Workshop, pp 87–98
- Jones SP et al. (2003) The Haskell 98 Language and Libraries: The Revised Report. J Func Prog 13(1):0–255. <http://www.haskell.org/definition/>
- Zhao L (2002) Solving and Creating Difficult Instances of Posts Correspondence Problem. Department of Computer Science, University of Alberta, Master's thesis
- Chakravarty M, Keller G, Jones SP (2005) Associated type synonyms. In: Proc. of the 10th ACM SIGPLAN International Conf. on Functional Programming (ICFP'05), pp 241–253
- Chakravarty M, Keller G, Jones SP, Marlow S (2005) Associated types withclass. In: Proc. of the ACM Symp. on Principles of Prog. Languages (POPL'05), pp 1–13
- Jones M (1994) Qualified Types. Cambridge University Press, Cambridge
- Jones M (1995) Simplifying and Improving Qualified Types. In: Proc. of the ACM Conf. on Functional Prog. and Comp. Architecture (FPCA'95), pp 160–169
- Jones M (2000) Type Classes with Functional Dependencies. In: Proc. of the European Symp. on Programming (ESOP'2000). LNCS 1782
- Sulzmann M, Duck G, Jones SP, Stuckey P (2007) Understanding functional dependencies via constraint handling rules. J Funct Program 17(1):83–129
- Milner R (1978) A theory of type polymorphism in programming. J Comput Syst Sci 17:348–375
- Stuckey P, Sulzmann M (2005) A Theory of Overloading. ACM Trans Prog Lang Syst (TOPLAS) 27(6):1216–1269
- Wadler P, Blott S (1989) How to make ad-hoc polymorphism less ad-hoc. In: Proc. of the 16th ACM Symp. on Principles of Prog. Lang. (POPL'89), pp 60–76. ACM Press, New York
- Jones SP et al (1998) GHC—The Glasgow Haskell Compiler. <http://www.haskell.org/ghc/>
- Jones SP et al (2011) GHC—The Glasgow Haskell Compiler 7.0.4 User's Manual. <http://www.haskell.org/ghc/>