

A Unification Algorithm for COQ Featuring Universe Polymorphism and Overloading

Beta Ziliani

MPI-SWS
beta@mpi-sws.org

Matthieu Sozeau

Inria & PPS, Université Paris Diderot
matthieu.sozeau@inria.fr

Abstract

Unification is a core component of every proof assistant or programming language featuring dependent types. In many cases, it must deal with higher-order problems up to conversion. Since unification in such conditions is undecidable, unification algorithms may include several heuristics to solve common problems. However, when the stack of heuristics grows large, the result and complexity of the algorithm can become unpredictable.

Our contributions are twofold: (1) We present a full description of a new unification algorithm for the Calculus of Inductive Constructions (the base logic of COQ), including universe polymorphism, canonical structures (the overloading mechanism baked into COQ’s unification), and a small set of useful heuristics. (2) We implemented our algorithm, and tested it on several libraries, providing evidence that the selected set of heuristics suffices for large developments.

Categories and Subject Descriptors D.1.1 [Programming Techniques]: Applicative (Functional) Programming; F.4.1 [Mathematical Logic And Formal Languages]: Mathematical Logic—Mechanical theorem proving

Keywords Interactive theorem proving; unification; Coq; universe polymorphism; overloading.

1. Introduction

In the last decade proof assistants have become more sophisticated and, as a consequence, increasingly adopted by computer scientists and mathematicians. In particular, they are being adopted to help dealing with very complex proofs, proofs that are hard to grasp—and more importantly, to *trust*—for a human. For example, in the area of algebra, the Feit-Thompson Theorem was recently formalized [10] in the proof assistant COQ [22]. To provide a sense of the accomplishment of Gonthier and his team, the original proof of this theorem was published in two volumes, totaling an astounding 250 pages. The team formalized it entirely in COQ, together with several books of algebra required as background material.

In order to make proofs manageable, this project relies heavily on the ability of COQ’s unification algorithm to infer implicit

arguments and expand heavily overloaded functions. This goes to the point that it is not rare to find in the source files a short definition that is expanded, by the unification algorithm, into several lines of code in the *Calculus of (co-)Inductive Constructions* (CIC), the base logic of COQ. This expansion is possible thanks to the use of the overloading mechanism in COQ called *canonical structures* [20]. This mechanism, similar in spirit to Haskell’s *type classes*, is baked into the unification algorithm. By being part of unification, this mechanism has a unique opportunity to drive unification to solve particular unification problems in a similar fashion to Matita’s *hints* [3]. It is so powerful, in fact, that it enables the development of dependently-typed logic meta-programs [12].

Another important aspect of the algorithm is that it must deal with higher-order problems, which are inherently undecidable, up to a subtyping relation on universes. For this reason, the current implementation of the unification algorithm has grown with several heuristics, yielding acceptable solutions to common problems in practice. Unfortunately, the algorithm is unpredictable and hard to reason about: given a unification problem, it is hard to predict the substitution the algorithm will return, and the time complexity for the task. This unpredictability of the current implemented algorithm has two main reasons: (i) it lacks a specification, and (ii) it incorporates a number of heuristics that obfuscate the order in which unification subproblems are considered.

While the algorithm being unpredictable is bad on its own, the problem gets exacerbated when combined with canonical structures, since their resolution may depend on the solutions obtained in previous unification problems. To somehow accommodate for this unfortunate situation, several works in the literature explain canonical structures by example [8, 9, 12, 13], providing some intuition on how canonical structures work, in some cases even detailing certain necessary aspects of the unification process. However, they fall short of explaining the complex process of unification as a whole.

This paper presents our remedy to the current situation. More precisely, our four main contributions are:

1. An original, full-fledged description of a unification algorithm for CIC, incorporating canonical structures and universe polymorphism [21].
2. The first formal description, to the best of our knowledge, of an extremely useful heuristic implemented in the unification algorithm of COQ, *controlled backtracking*.
3. A corresponding pluggable implementation, incorporating only a restricted set of heuristics, such as controlled backtracking. Most notably, we purposely left out a technique known as *constraint postponement*, present in many systems and in the current implementation in Coq, which may reorder unification subproblems. This reordering prevents us from knowing exactly when equations are being solved.

$\text{in_head} : \forall (x : A) (l : \text{list } A), x \in (x :: l)$
 $\text{in_tail} : \forall (x : A) (y : A) (l : \text{list } A), x \in l \rightarrow x \in (y :: l)$

Lemma $\text{inL} : \forall (x : A) (l r : \text{list } A), x \in l \rightarrow x \in (l ++ r)$
Lemma $\text{inR} : \forall (x : A) (l r : \text{list } A), x \in r \rightarrow x \in (l ++ r)$

Figure 1. List membership axioms and lemmas.

- Evidence that such principled heuristics suffice to solve 99.9% of the unification problems that arise in libraries such as the Mathematical Components library [11] and CPDT [7].

It is interesting to note that during this work we found two bugs in the logic of the original unification algorithm of COQ. While this work focuses on the COQ proof assistant, the problems and solutions presented may be of interest to other type theory based assistants and programming languages, such as Agda [16], Matita [2], or Idris [5].

In the rest of the paper, we start introducing with examples some features and heuristics included in COQ’s unification algorithm (§2). Then, we present the language used in the paper (§3), necessary to understand the core contribution of this work, a new unification algorithm (§4). We evaluate the algorithm (§5) and conclude (§6).

2. COQ’s Unification at a Glance

We start by showing little examples highlighting some of the particularities of COQ’s unification algorithm.

First-order approximation: In many cases, a unification problem may have several incomparable solutions. Consider for instance the following definition in a context where y_1 and y_2 are defined:

Definition $\text{ex0} : y_1 \in ([y_1] ++ [y_2]) := \text{inL} _ _ _ (\text{in.head} _ _)$

We assume the definitions and lemmas for list membership listed in Figure 1, and note $(x :: s)$ for the *consing* of x to list s , $[]$ for the empty list, and $l ++ r$ for the concatenation of lists l and r . We also denote $[a_1; \dots; a_n]$ a list with elements a_1 to a_n .

This definition is a proof that the element y_1 is in the list resulting from concatenating the singleton lists $[y_1]$ and $[y_2]$. The proof in itself provides evidence that the element is in the head (in.head) of the list on the left (inL). As customary in COQ code, the type annotation shows what the definition is proving, and the proof omits the information that can be inferred, replacing each argument to inL and in.head with *holes* ($_$). The elaboration mechanism of COQ, that is, the algorithm in charge of filling up these holes, calls the unification algorithm with the following unification problem, where the left-hand side corresponds to what the body of the definition proves, and the right-hand side to what it is expected to prove:¹

$$?z_1 \in ((?z_1 :: ?z_2) ++ ?z_3) \approx y_1 \in ([y_1] ++ [y_2])$$

where $?z_1, ?z_2$ and $?z_3$ are fresh meta-variables. In turn, after assigning y_1 to $?z_1$, the unification algorithm has to solve the following problem:

$$(y_1 :: ?z_2) ++ ?z_3 \approx [y_1] ++ [y_2]$$

One possible solution to this equation is to assign $[]$ to $?z_2$, and $[y_2]$ to $?z_3$, which corresponds to equate each argument of the concatenation, similar to what we did before with the \in *predicate*. However, since concatenation is a *function*, *i.e.*, it computes the

¹How elaboration works will not be discussed in this work. The interested reader is invited to read [4], which provides details on bi-directional elaboration in the Matita proof assistant, also based on CIC.

concatenation of the two lists, there are other possible solutions that makes both terms *convertible* (*i.e.*, having the same normal form). One such solution, for instance, is to assign $[y_2]$ to $?z_2$, and $[]$ to $?z_3$.

Many works in the literature [*e.g.*, 1, 14, 17, 18] are devoted to the creation of unification algorithms returning a *Most General Unifier* (MGU), that is, a *unique* solution that serves as a representative for all *convertible* solutions. Agda [16], for instance, which incorporates such type of unification algorithm, fails to compile Example `ex0` above, since no such MGU exists. This forces the proof developer to manually fill-in the holes.

Despite the equation having multiple solutions, however, not every solution is equally “good”. For `ex0`, the first solution is the most *natural* one, meaning the one expected by the proof developer. For this reason, instead of failing, COQ favors syntactic equality by trying first-order unification. Formally, when faced with a problem of the form

$$t \ t_1 \ \dots \ t_n \ \approx \ u \ u_1 \ \dots \ u_n$$

the algorithm decomposes the problem into $n + 1$ subproblems, first equating $t \approx u$, and then $t_i \approx u_i$, for $0 < i \leq n$.

Controlled backtracking: In [19, chp. 10], a unification algorithm for CIC is presented, performing *only* first-order unification. In COQ, instead, when first-order approximation fails, in an effort to find a solution to the equation, the algorithm reduces the terms *carefully*. For instance, consider the following variation of the previous example, where the list on the left of the concatenation is **let-bound**:

Definition $\text{ex1} : y_1 \in (\text{let } l := [y_1] \text{ in } (l ++ [y_2]))$
 $:= \text{inL} _ _ _ (\text{in.head} _ _)$

The main equation to solve now is

$$(y_1 :: ?z_2) ++ ?z_3 \approx \text{let } l := [y_1] \text{ in } (l ++ [y_2])$$

Since both terms do not share the same *head* (the concatenation operator on the left and the **let**-binding on the right), the algorithm reduces the **let**-binding, obtaining the same problem as in `ex0`. Note that it has to be careful: it should not reduce the concatenation operator, otherwise the problem will become unsolvable. For this reason, it delays the unfolding of constants, such as $++$, and, in the case of having constants on both sides of the equation, it takes special care of which one to unfold. This heuristic enables fine control over the instance resolution mechanism of canonical structures [12].

Canonical structures: *Canonical structures* (CS) is a powerful overloading mechanism, baked into the unification algorithm. We demonstrate this mechanism with a typical example from overloading: the equality operator. Similar to how type classes are used in Haskell [23], we define a *class* or, in CS terminology, a *structure*:²

Structure $\text{eqType} := \text{EqType} \{ \text{sort} : \text{Type};$
 $\text{equal} : \text{sort} \rightarrow \text{sort} \rightarrow \text{bool} \}$

eqType is a record type with two fields: a type sort , and a boolean binary operation equal on sort . These fields can be accessed using projectors:

$\text{sort} : \text{eqType} \rightarrow \text{Type}$
 $\text{equal} : \forall e : \text{eqType}. \text{sort } e \rightarrow \text{sort } e \rightarrow \text{bool}$

To construct an element of the type, the constructor EqType is provided, which takes the values for the two fields as arguments. For example, one possible eqType *instance* for bool is:

Definition $\text{eqType.bool} := \text{EqType } \text{bool } \text{eq.bool}$

²This example is a significant simplification of one taken from [11, 12].

where $\text{eq_bool } x \ y := (x \ \&\& \ y) \ || \ (!x \ \&\& \ !y)$. (We denote boolean conjunction, disjunction and negation as $\&\&$, $||$ and $!$.)

Similarly, it is possible to declare *recursive* instances. For example, consider the instance for the pair type $A \times B$, where A and B are themselves instances of eqType :

Definition $\text{eqType_pair } (A \ B : \text{eqType}) :=$
 $\text{EqType } (\text{sort } A \times \text{sort } B) \ (\text{eq_pair } A \ B)$

where

$\text{eq_pair } (A \ B : \text{eqType}) \ (u \ v : \text{sort } A \times \text{sort } B) :=$
 $(\text{equal } A \ (\pi_1 \ u) \ (\pi_1 \ v)) \ \&\& \ (\text{equal } B \ (\pi_2 \ u) \ (\pi_2 \ v))$

In order to use instances eq_bool and eq_pair for overloading, we need to declare them as **Canonical**. After they have been declared canonical, whenever the elaboration mechanism is asked to elaborate a term like $\text{equal } _ \ (b_1, b_2) \ (c_1, c_2)$, for booleans b_1, b_2, c_1 and c_2 , it will generate a unification problem matching the expected and inferred type of the second argument of equal , that is,

$\text{sort } ?e \approx \text{bool} \times \text{bool}$

for some meta-variable $?e$ elaborated from the *hole* $(_)$.

To solve the equation above, COQ's unification will try instantiating $?e$ using the canonical instance eqType_pair , resulting in two new unification subproblems, for fresh meta-variables $?A$ and $?B$:

$\text{sort } ?A \approx \text{bool} \quad \text{sort } ?B \approx \text{bool}$

Next, it will choose $?A := \text{eqType_bool}$ and $?B := \text{eqType_bool}$, resulting in that $\text{equal } ?e \ (b_1, b_2) \ (c_1, c_2)$ reduces, as expected, to $\text{eq_bool } b_1 \ c_1 \ \&\& \ \text{eq_bool } b_2 \ c_2$.

We can declare a number of canonical eqType instances for our types equipped with decidable equality. Then, we can uniformly write $\text{equal } _ \ t \ u$, and let unification compute the corresponding instance for the hole, according to the type of t and u .

Polymorphic universes and subtyping: Unification in CIC is not a simple equational theory, in the sense that it must deal with the subtyping relation generated by the cumulative universe hierarchy ($\text{Type}(i) \leq \text{Type}(j) \iff i \leq j$). To our knowledge, we present the first algorithm dealing with this relation properly. In COQ, previous algorithms relied on the kernel to check the proper use of universes, resulting in particular in non-local error reporting and the inability to backtrack on these errors, which becomes crucial in presence of universe polymorphism and first-order approximation.

3. The Language: CIC with Open Terms

Before presenting the algorithm, we need to present the base language of COQ, the Calculus of Inductive Constructions (CIC) [22, chap. 4]. It is a dependently typed λ -calculus extended with inductive types. It also includes co-inductive types, but their formulation is not important for this work, so it will be omitted.

The terms (and types) of the language are defined as

$t, u, T, U = x \ | \ c[\bar{\ell}] \ | \ i[\bar{\ell}] \ | \ k[\bar{\ell}] \ | \ s \ | \ ?x[\sigma]$
 $\quad | \ \forall x : T. U \ | \ \lambda x : T. t \ | \ t \ u \ | \ \mathbf{let} \ x := t : T \ \mathbf{in} \ u$
 $\quad | \ \mathbf{match}_T \ t \ \mathbf{with} \ k_1 \ \bar{x}_1 \Rightarrow t_1 \ | \ \dots \ | \ k_n \ \bar{x}_n \Rightarrow t_n \ \mathbf{end}$
 $\quad | \ \mathbf{fix}_j \ \{x_1/n_1 : T_1 := t_1; \dots; x_m/n_m : T_m := t_m\}$
 $\sigma = \bar{\ell}$
 $\ell, \kappa \in \mathcal{L} \cup 0^-$
 $K = \kappa \ | \ K + 1$
 $s = \text{Type}(\bar{K}^+)$

Terms include variables $x \in \mathcal{V}$, constants $c \in \mathcal{C}$, inductive type constructors $i \in \mathcal{I}$ and constructors $k \in \mathcal{K}$, these last three being

applied to universe instances $\bar{\ell}$ built from universe levels $\ell \in \mathcal{L} \cup 0^-$. Terms also includes sorts s , representing algebraic universes. Algebraic universes represent least upper bounds of a (non-empty) set of levels or successors of levels. They are used notably to sort products, e.g. $(\forall A : \text{Type}(i), \text{Type}(j)) : \text{Type}(i+1, j)$. The impredicative sort Prop , the type of propositions, is represented as $\text{Type}(0^-)$. Terms may contain a *hole*, representing a missing piece of the term (or proof). Holes are represented with meta-variables, a variable prepended with a question mark, as in $?x$. For reasons that will become apparent soon, meta-variables are applied to a *suspended substitution* $[\sigma]$, which is nothing more than a list of terms.

In order to destruct an element of an inductive type, CIC provides regular pattern **matching** and mutually recursive **fixpoints**. Their notation is slightly different from, but easily related to, the actual notation from COQ. **match** is annotated with the return predicate T , meaning that the type of the whole **match** expression may depend on the element being pattern matched (**as** ... **in** ... in standard COQ notation). In the **fix** expression, $x/n : T := t$ means that T is a type starting with at least n product types, and the n -th variable is the decreasing one in t (**struct** in COQ notation). The subscript j of **fix** selects the j -th function as the main entry point of the mutually recursive fixpoints.

In order to typecheck and reduce terms, COQ uses several contexts, each handling different types of knowledge:

1. Universe contexts Φ declaring universe level names and associated constraints ((in-)equalities on levels);
2. Local contexts Γ , including bound variables and **let**-bound expressions;
3. Meta-contexts Σ , containing meta-variable declarations and definitions; and
4. A global environment E , containing the global knowledge; that is, axioms, theorems, and inductive definitions, along with a *global* universe context that can be incrementally enriched.

Formally, they are defined as follows:

$\Phi = \bar{\ell} \models \mathcal{C} \quad \mathcal{C} = \cdot \ | \ \mathcal{C} \wedge \ell \ \mathcal{O} \ \ell' \quad \text{where } \mathcal{O} \in \{=, \leq, <\}$
 $\Gamma, \Psi = \cdot \ | \ x : T, \Gamma \ | \ x := t : T, \Gamma$
 $\Sigma = \cdot \ | \ ?x : T[\Psi], \Sigma \ | \ ?x := t : T[\Psi], \Sigma$
 $E = \cdot \ | \ c[\Phi] : T, E \ | \ c[\Phi] := t : T, E \ | \ I, E \ | \ \Phi, E$
 $I = \forall \Phi, \Gamma. \{ i : \forall \bar{y} : \bar{T}_h. s := \{k_1 : U_1; \dots; k_n : U_n\} \}$

A universe context consists of a list of levels $\bar{\ell}$ and a set of constraints \mathcal{C} on those levels. The local context is standard, and requires no further explanation. Meta-variables have *contextual types*, meaning that the type T of a meta-variable must have all of its free variables bound within the local context Ψ . In this work we borrow the notation $T[\Psi]$ from Contextual Modal Type Theory [15]. A meta-variable can be instantiated with a term t , noted $?x := t : T[\Psi]$. In this case, t should also contain only free variables occurring in Ψ .

The global environment associates a constant c with a local universe context (usually omitted), a type and, optionally, a definition. In the first case, c is an axiom, while in the second c is a theorem proved by term t . Additionally, this environment may also contain (mutually recursive) inductive types and global universe and constraint declarations.

A set of mutually recursive inductive types I is prepended with a universe context Φ and a list of parameters Γ . Every inductive type i defined in the set has sort s , with parameters $\bar{y} : \bar{T}_h$. It has a possibly empty list of constructors k_1, \dots, k_n . For every j , each type U_j of constructor k_j has shape $\forall \bar{z} : \bar{U}'. i \ t_1 \ \dots \ t_n$.

Inductive definitions are restricted to avoid circularity, meaning that every type constructor i can only appear in a strictly positive position in the type of every constructor. For the purpose of this work, understanding this restriction is not crucial, and we refer the interested reader to [22, chap. 4]. Additionally, fixpoints on inductive types must pass the guard condition (ibid., §4.5.5) to be accepted by the kernel, a syntactic criterion ensuring termination. We will come back to this point in §4.7. It is interesting to note that a **Structure**, like the one shown in the previous section, is syntactic sugar for an inductive type with only one constructor, and with projections generated for each argument of the constructor.

3.1 Meta-Variables and Contextual Types

At a high-level, meta-variables are holes in a term, which are expected to be filled out at a later point in time. For instance, when a lemma is applied to solve some goal, COQ internally creates fresh meta-variables for all the formal parameters of the lemma, and proceeds to unify the goal with the conclusion of the lemma. During unification, meta-variables are instantiated so that both terms (the goal and the conclusion of the lemma) become *convertible* (equal modulo reduction rules, see §3.2).

In the simple examples shown so far, contextual types played no role but, as we are going to see in the next example, they prevent illegal instantiations of meta-variables. For instance, such illegal instantiations could potentially happen if the same meta-variable occurs at different locations in a term, with different variables in the scope of each occurrence. We illustrate this point with an example taken from [24]. Suppose function f defined locally as follows:

$$f := \lambda w : \text{nat}. (_ : \text{nat})$$

where the hole $(_)$ is an indication to COQ's elaboration mechanism to "fill in this hole with a meta-variable". The accessory typing annotation provides the expected type for the meta-variable. Assuming no other variables occur in scope, after elaboration f becomes:

$$f := \lambda w : \text{nat}. ?v[w] \quad (1)$$

for some fresh meta-variable $?v$. Since any instantiation of $?v$ may only refer to w , its type becomes $\text{nat}[w : \text{nat}]$. This contextual type specifies precisely that $?v$ may only be instantiated with a term of type nat containing at most a single free variable w of type nat . In the elaborated term (1), $[w]$ stands for the *suspended substitution* specifying how to transform such instantiation into one that is well-typed under the current context. In this case, this substitution is the identity, because the current context and the context under which $?v$ was created are identical (in fact, the latter is a copy of the former).

Now suppose that we define functions g and h referring to f :

$$g := \lambda x y : \text{nat}. f x \quad h := \lambda z : \text{nat}. f z$$

and proceed to unify g with a function projecting the first argument:

$$g \approx \lambda x y : \text{nat}. x$$

In order to solve this equation, COQ proceeds to unfold the definition of g and to push x and y in the local context. The new equation to solve becomes:

$$f x \approx x$$

After unfolding f and β -reducing the left-hand side, it amounts to solving the following equation:

$$?v[x] \approx x$$

At this point is where the contextual type of $?v$ comes into play. If meta-variables were created with a normal type, that is, not having contextual type (and suspended substitution), it would seem that the only solution for $?v$ is x . However, that solution would break the definition of h since x is not in scope there. Given the contextual information, however, COQ will correctly realize that

$$(\lambda x : T. t) u \rightsquigarrow_{\beta} t\{u/x\} \quad \mathbf{let} x := u : T \mathbf{in} t \rightsquigarrow_{\zeta} t\{u/x\}$$

$$\frac{(x := t : T) \in \Gamma}{x \rightsquigarrow_{\delta\Gamma} t} \quad \frac{?x := t : T[\Psi] \in \Sigma}{?x[\sigma] \rightsquigarrow_{\delta\Sigma} t\{\sigma/\widehat{\Psi}\}}$$

$$\frac{(c[\bar{\ell}] \vDash \mathcal{C}) := t : T) \in E}{c[\bar{\kappa}] \rightsquigarrow_{\delta E} t[\bar{\kappa}/\bar{\ell}]}$$

$$\mathbf{match}_T k_j[\bar{\kappa}] \bar{t} \mathbf{with} \bar{k} \bar{x} \Rightarrow \bar{u} \mathbf{end} \rightsquigarrow_{\iota} u_j\{\bar{t}/x_j\}$$

$$\frac{F = x/n : T := t \quad a_n = k_j[\bar{\kappa}] \bar{t}}{\mathbf{fix}_j \{F\} \bar{a} \rightsquigarrow_{\iota} t_j\{\mathbf{fix}_m \{F\}/x_m\} \bar{a}}$$

Figure 2. Reduction rules in CIC.

$?v$ should be instantiated with w , not x . Under that instantiation, g will normalize to $\lambda x y : \text{nat}. x$, and h will normalize to $\lambda z : \text{nat}. z$.

The suspended substitution and the contextual type are the tools that the unification algorithm uses to know how to instantiate the meta-variable. The decision to solve $?v[x] \approx x$ by instantiating $?v : \text{nat}[w : \text{nat}]$ with w is due to the problem falling in the *pattern* unification subset [14]. When COQ faces a problem of the form

$$?u[y_1, \dots, y_n] \approx e$$

where the y_1, \dots, y_n are all distinct variables, then the *most general* solution to the problem is to *invert* the substitution and apply it on the right-hand side of the equation, in other words instantiating $?u$ with $e\{x_1/y_1, \dots, x_n/y_n\}$, where x_1, \dots, x_n are the variables in the local context of $?u$ (and assuming the free variables of e are in $\{y_1, \dots, y_n\}$).

In the example above, at the point where COQ tries to unify $?u[x] \approx x$, the solution (through inversion) is to instantiate $?u$ with $x\{w/x\}$, that is, w .

3.2 Semantics

Reduction of CIC terms is performed through a set of rules listed in Figure 2. Besides the standard β rule, CIC provides six more rules to destruct the different term constructions: the ζ rule, which expands let-definitions, three δ rules, which expand definitions from each of the contexts, and two ι rules, which evaluate pattern **matchings** and **fixpoints**.

Most of the rules are self explanatory, with the sole exception of the $\delta\Sigma$ rule. It takes a meta-variable $?x$, applied to suspended substitution σ , and replaces it by its definition t , replacing each variable from its local context Ψ by the corresponding term from substitution σ . For this we use the multi-substitution of terms, mapping the variables coming from the domain of Ψ with terms in σ . To obtain the domain of Ψ , we use the type-eraser function $\widehat{\cdot}$, defined as:

$$x_1 : T_1, \dots, x_n : T_n = x_1, \dots, x_n$$

The *unfolding* rules ($\delta\Gamma$, $\delta\Sigma$, δE), of course, depend on the contexts. As customary, we will always consider the environment E implicit. We will also omit Γ and Σ when there is no room for ambiguity.

Conversion (\equiv) is defined as the congruent closure of these reduction rules, plus η -conversion: $u \equiv \lambda x : T. u x$ iff $x \notin \text{FV}(u)$.

4. The Algorithm

We proceed to describe our proposal in the following pages, underlining every non-standard design decision. We emphasize that this

is not a faithful description of the current unification algorithm in COQ, but a new one that is, however, close enough. In particular, we purposely left out a technique known as *constraint postponement* (§4.6), as well as other heuristics hard to grasp. At the same time, we introduced our own set of heuristics, based on practical examples (§5) and on previous work by Abel and Pientka (§4.3).

The unification judgment is of the form:

$$\Phi; \Sigma; \Gamma \vdash t_1 \approx_{\mathcal{R}} t_2 \triangleright \Phi', \Sigma'$$

It unifies terms t_1 and t_2 , given a universe context Φ , meta-context Σ and a local context Γ . There is an implicit global environment E . The universe context carries additional information for each universe variable introduced: they are either flexible (ℓ_x) or rigid (ℓ_r). This information is used when unifying two instances of the same constant to avoid forcing universe constraints that would not appear if the bodies of the instantiations were unified instead, respecting transparency of the constants. Flexible variables are generated when taking a fresh instance of a polymorphic constant, inductive or constructor during elaboration, while rigid ones correspond to user-specified levels or Type annotations. The relation \mathcal{R} (\equiv or \leq) indicates if we are trying to derive conversion of the two terms or *cumulativity*, the subtyping relation on universes. The rules decomposing constructions switch to conversion in their premises except for sorts and dependent function spaces, otherwise the relation is preserved when reducing one side or the other. The algorithm returns a new universe context Φ' and meta-context Σ' , which are extensions of Φ and Σ , respectively, perhaps with new universes constraints in Φ' , and new meta-variables or instantiations of existing meta-variables in Σ . The algorithm ensures that terms t_1 and t_2 are convertible (or in the cumulativity relation) in the returned contexts.

In the presentation of the algorithm we will often omit the universe context Φ , since it is, in most of the cases, simply threaded along. The unification algorithm is quite involved, so to help readability we split the rules across four different subsections. Roughly, in §4.1 we consider the case when the two terms being unified have no arguments, and share the same head constructor; in §4.2 we consider terms having arguments; in §4.3 we consider meta-variable unification; and in §4.4 we consider canonical structures resolution. The algorithm's strategy, which backtracks in some particular cases, cannot be understood by the ordering of the rules, so we devote §4.5 to explain in detail the algorithm's strategy. In §4.6 we explain the technique of constraint postponement, and the reason for its omission in our algorithm. Finally, in §4.7 we comment on the correctness of the algorithm.

4.1 Same Constructor

Figure 3 shows the rules that apply when both terms share the same head constructor. We need to distinguish this set of rules from the other rules in the algorithm, so we annotate them with a 0 as subscript of the turnstile (\vdash_0). The reasons will become evident when we look at the rules in the next subsection.

The rule **TYPE-SAME** unifies two sorts, according to the relation \mathcal{R} . By invariant, we know that the right-hand side universe can only be a single level while the l.h.s. can be the least upper bound of a set of universe levels or successors iff the relation is cumulativity, and any such \leq constraints can be translated to a set of atomic \leq or $<$ constraints (see [21] for details). The predicate $\mathcal{C} \models$ denotes satisfiability of set of constraints \mathcal{C} and naturally extends to consistency of universe contexts ($\phi \models$).

For abstractions (**LAM-SAME**) and products (**PROD-SAME**), we first unify the types of the arguments, and then the body of the binder, with the local context extended with the bound variable. (The universe context is omitted, as in some of the following rules, since it is just threaded along.) When unifying two **lets**, the rule

$$\text{TYPE-SAME} \quad \frac{\mathcal{C}' = \mathcal{C} \wedge \bar{u} \mathcal{R} \kappa \quad \mathcal{C}' \models}{\ell \models \mathcal{C}; \Sigma; \Gamma \vdash_0 \text{Type}(\bar{u}) \approx_{\mathcal{R}} \text{Type}(\kappa) \triangleright \ell \models \mathcal{C}'; \Sigma}$$

PROD-SAME, LAM-SAME

$$\frac{\Pi \in \{\lambda, \forall\} \quad \Sigma_0; \Gamma \vdash T_1 \approx_{\equiv} U_1 \triangleright \Sigma_1 \quad \Sigma_1; \Gamma, x : T_1 \vdash T_2 \approx_{\mathcal{R}} U_2 \triangleright \Sigma_2}{\Sigma_0; \Gamma \vdash_0 \Pi x : T_1. T_2 \approx_{\mathcal{R}} \Pi x : U_1. U_2 \triangleright \Sigma_2}$$

LET-SAME

$$\frac{\Sigma_0; \Gamma \vdash T \approx_{\equiv} U \triangleright \Sigma_1 \quad \Sigma_1; \Gamma \vdash t_2 \approx_{\equiv} u_2 \triangleright \Sigma_2 \quad \Sigma_2; \Gamma, x := t_2 \vdash t_1 \approx_{\mathcal{R}} u_1 \triangleright \Sigma_3}{\Sigma_0; \Gamma \vdash_0 \text{let } x := t_2 : T \text{ in } t_1 \approx_{\mathcal{R}} \text{let } x := u_2 : U \text{ in } u_1 \triangleright \Sigma_3}$$

RIGID-SAME

$$\frac{h \in \mathcal{V} \cup \mathcal{I} \cup \mathcal{K} \quad C_1 = C_0 \wedge \bar{\kappa} = \bar{\kappa}' \quad C_1 \models}{(\bar{\ell} \models C_0); \Sigma; \Gamma \vdash_0 h[\bar{\ell}] \approx_{\mathcal{R}} h[\bar{\kappa}'] \triangleright (\bar{\ell} \models C_1), \Sigma}$$

FLEXIBLE-SAME

$$\frac{h \in \mathcal{C} \quad \Phi_0 \models \bar{\ell} = \bar{\kappa} \triangleright \Phi_1}{\Phi_0; \Sigma; \Gamma \vdash_0 h[\bar{\ell}] \approx_{\mathcal{R}} h[\bar{\kappa}] \triangleright \Phi_1, \Sigma}$$

UNIV-EQ

$$\frac{\Phi \models i = j}{\Phi \models i = j \triangleright \Phi}$$

UNIV-FLEXIBLE

$$\frac{i_x \vee j_x \in \bar{\ell} \quad C \wedge i = j \models}{(\bar{\ell} \models C) \models i = j \triangleright (\ell \models C \wedge i = j)}$$

CASE-SAME

$$\frac{\Sigma_0; \Gamma \vdash T \approx_{\equiv} U \triangleright \Sigma_1 \quad \Sigma_1; \Gamma \vdash t \approx_{\equiv} u \triangleright \Sigma_2 \quad \Sigma_2; \Gamma \vdash \bar{b} \approx_{\equiv} \bar{b}' \triangleright \Sigma_3}{\Sigma_0; \Gamma \vdash_0 \text{match}_T t \text{ with } \bar{b} \text{ end} \approx_{\mathcal{R}} \text{match}_U u \text{ with } \bar{b}' \text{ end} \triangleright \Sigma_3}$$

FIX-SAME

$$\frac{\Sigma_0; \Gamma \vdash \bar{T} \approx_{\equiv} \bar{U} \triangleright \Sigma_1 \quad \Sigma_1; \Gamma \vdash \bar{t} \approx_{\equiv} \bar{u} \triangleright \Sigma_2}{\Sigma_0; \Gamma \vdash_0 \text{fix}_j \{x/n : T := t\} \approx_{\mathcal{R}} \text{fix}_j \{x/n : U := u\} \triangleright \Sigma_2}$$

Figure 3. Unifying terms sharing the same head constructor.

LET-SAME compares first the type of the definitions, then the definitions themselves, and finally the body. In the last case, it augments the local context with the definition on the left (choosing the one on the left is an arbitrary choice, but after unification both definitions are convertible, *i.e.*, indistinguishable).

RIGID-SAME equates the same variable, inductive type or constructor, enforcing that their universe instances are equal (note that the application of the rule will fail if these new constraints are inconsistent). The **FLEXIBLE-SAME** rule unifies two instances of the same constant using a stronger condition on universe instances: they must unify according to the current constraints and by equating rigid universe variables with flexible variables only ($\Phi \models i = j$ checks if the constraint is already derivable). Otherwise we will backtrack on this rule to unfold the constant and unify the bodies (§4.2), which will generally result in weaker, more general constraints to be enforced. The last two rules (**CASE-SAME** and **FIX-SAME**) unify **matches** and **fixpoints**, respectively. In both cases we just unify pointwise every component of the term constructors.

4.2 Reduction

The previous subsection considered only the cases when both terms have no arguments and share the same constructor. If that is not the case, the algorithm first tries first-order approximation (rule **APP-FO** in Figure 4). This rule, when considering two applications with the same number of arguments (n), compares the head element

$$\text{APP-FO} \quad \frac{\Sigma_0; \Gamma \vdash_0 t \approx_{\mathcal{R}} u \triangleright \Sigma_1 \quad n \geq 0 \quad \Sigma_1; \Gamma \vdash \bar{t}_n \approx_{\equiv} \bar{u}_n \triangleright \Sigma_2}{\Sigma_0; \Gamma \vdash t \bar{t}_n \approx_{\mathcal{R}} u \bar{u}_n \triangleright \Sigma_2}$$

$$\begin{array}{c} \text{META-}\delta\text{R, LAM-}\beta\text{R, LET-}\zeta\text{R} \\ \Sigma; \Gamma \vdash u \overset{w}{\rightsquigarrow}_{\delta\Sigma, \beta, \zeta} u' \\ \Sigma; \Gamma \vdash t \approx_{\mathcal{R}} u' \triangleright \Sigma' \\ \hline \Sigma; \Gamma \vdash t \approx_{\mathcal{R}} u \triangleright \Sigma' \end{array} \quad \begin{array}{c} \text{META-}\delta\text{L, LAM-}\beta\text{L, LET-}\zeta\text{L} \\ \Sigma; \Gamma \vdash t \overset{w}{\rightsquigarrow}_{\delta\Sigma, \beta, \zeta} t' \\ \Sigma; \Gamma \vdash t' \approx_{\mathcal{R}} u \triangleright \Sigma' \\ \hline \Sigma; \Gamma \vdash t \approx_{\mathcal{R}} u \triangleright \Sigma' \end{array}$$

$$\text{CASE-}\iota\text{R} \quad \frac{u \text{ is } \mathbf{fix} \text{ or } \mathbf{match} \quad \Sigma; \Gamma \vdash u \downarrow_{\beta\zeta\delta\Sigma\iota\theta}^w u' \quad u \neq u' \quad \Sigma; \Gamma \vdash t \approx_{\mathcal{R}} u' \triangleright \Sigma'}{\Sigma; \Gamma \vdash t \approx_{\mathcal{R}} u \triangleright \Sigma'}$$

$$\text{CASE-}\iota\text{L} \quad \frac{t \text{ is } \mathbf{fix} \text{ or } \mathbf{match} \quad \Sigma; \Gamma \vdash t \downarrow_{\beta\zeta\delta\Sigma\iota\theta}^w t' \quad t \neq t' \quad \Sigma; \Gamma \vdash t' \approx_{\mathcal{R}} u \triangleright \Sigma'}{\Sigma; \Gamma \vdash t \approx_{\mathcal{R}} u \triangleright \Sigma'}$$

$$\begin{array}{c} \text{CONS-}\delta\text{NOTSTUCKR} \\ \text{not } \Sigma; \Gamma \vdash \text{is_stuck } u \\ u \overset{w}{\rightsquigarrow}_{\delta E, \delta\Gamma} u' \\ \Sigma; \Gamma \vdash t \approx_{\mathcal{R}} u' \triangleright \Sigma' \\ \hline \Sigma; \Gamma \vdash t \approx_{\mathcal{R}} u \triangleright \Sigma' \end{array} \quad \begin{array}{c} \text{CONS-}\delta\text{STUCKL} \\ \Sigma; \Gamma \vdash \text{is_stuck } u \\ t \overset{w}{\rightsquigarrow}_{\delta E, \delta\Gamma} t' \\ \Sigma; \Gamma \vdash t' \approx_{\mathcal{R}} u \triangleright \Sigma' \\ \hline \Sigma; \Gamma \vdash t \approx_{\mathcal{R}} u \triangleright \Sigma' \end{array}$$

$$\begin{array}{c} \text{CONS-}\delta\text{R} \\ \Sigma; \Gamma \vdash u \overset{w}{\rightsquigarrow}_{\delta E, \delta\Gamma} u' \\ \Sigma; \Gamma \vdash t \approx_{\mathcal{R}} u' \triangleright \Sigma' \\ \hline \Sigma; \Gamma \vdash t \approx_{\mathcal{R}} u \triangleright \Sigma' \end{array} \quad \begin{array}{c} \text{CONS-}\delta\text{L} \\ \Sigma; \Gamma \vdash t \overset{w}{\rightsquigarrow}_{\delta E, \delta\Gamma} t' \\ \Sigma; \Gamma \vdash t' \approx_{\mathcal{R}} u \triangleright \Sigma' \\ \hline \Sigma; \Gamma \vdash t \approx_{\mathcal{R}} u \triangleright \Sigma' \end{array}$$

$$\text{LAM-}\eta\text{R} \quad \frac{\begin{array}{c} u \text{'s head is not an abstraction} \\ \Sigma_0; \Gamma \vdash u : U \quad \text{ensure_product}(\phi_0; \Sigma_0; \Gamma; T; U) = (\phi_1; \Sigma_1) \\ \phi_1; \Sigma_1; \Gamma, x : T \vdash u x \approx_{\equiv} t \triangleright \phi_2; \Sigma_2 \end{array}}{\phi_0; \Sigma_0; \Gamma \vdash u \approx_{\mathcal{R}} \lambda x : T. t \triangleright \phi_2; \Sigma_2}$$

$$\text{LAM-}\eta\text{L} \quad \frac{\begin{array}{c} u \text{'s head is not an abstraction} \\ \Sigma_0; \Gamma \vdash u : U \quad \text{ensure_product}(\phi_0; \Sigma_0; \Gamma; T; U) = (\phi_1; \Sigma_1) \\ \phi_1; \Sigma_1; \Gamma, x : T \vdash t \approx_{\equiv} u x \triangleright \phi_2; \Sigma_2 \end{array}}{\phi_0; \Sigma_0; \Gamma \vdash \lambda x : T. t \approx_{\mathcal{R}} u \triangleright \phi_2; \Sigma_2}$$

Figure 4. Reduction steps attempted during unification.

(t and t' , using only the rules in Figure 3), and then proceeds to unify each of the arguments. As customary, we denote multiple applications as a *spine* [6], using the form $t \bar{u}_n$ to represent the term $(\dots (t u_1) \dots u_n)$. We call t the *head* of the term.

If the rules in Figure 3 plus APP-FO fail to apply, then the algorithm tries different reduction strategies. Except in some particular cases, the algorithm first tries reducing the right-hand side (rules ending with R) and then the left-hand side (rules ending with L). Except where noted, every L rule is just the mirror of the corresponding R rule, swapping the terms being unified in the conclusion and applications of unification in the premises. We will often omit the last letter (R or L), and simply write *e.g.*, META- δ when referring to both rules.

The algorithm first tries one step of either weak-head expansion of meta-variables ($\delta\Sigma$), weak-head β reduction, or weak-head **let**-expansion (ζ). These steps are described in rules META- δ , LAM- β ,

$$\frac{t \downarrow_{\beta\zeta\delta\iota}^w k_j \bar{a}}{\text{match}_T t \text{ with } \bar{k} \bar{x} \Rightarrow t' \text{ end}} \rightsquigarrow_{\theta} \text{match}_T k_j [\bar{k}] \bar{a} \text{ with } k \bar{x} \Rightarrow t' \text{ end}$$

$$\frac{a_{n_j} \downarrow_{\beta\zeta\delta\iota}^w k \bar{b}}{\mathbf{fix}_j \{F\} a_1 \dots a_{n_j} \rightsquigarrow_{\theta} \mathbf{fix}_j \{F\} a_1 \dots a_{n_j-1} (k \bar{b})}$$

Figure 5. The θ -reduction strategy.

and LET- ζ . (Actually, as we are going to see in §4.5, the order of the rules is slightly different, although for the moment the implicit ordering obtained from the figure suffices.)

More interesting are the cases for δE , $\delta\Gamma$ and ι reductions. The high level idea is that special care should be taken when unfolding defined constants and variables. One reason is efficiency: we hope that, before performing the unfolding of a constant or variable, we will find the same constant or variable on the other side of the equation. The second reason is to avoid missing potential solutions, as already mentioned when introducing controlled backtracking in §2.

In the case of a **match** or a **fix** (rules CASE- ι), we want to be able to reduce the scrutinee using all reduction rules, including δE and $\delta\Gamma$, and then (if applicable), continue reducing the corresponding branch of the **match** or the body of the **fix**, but avoiding δE and $\delta\Gamma$. We illustrate this desired behavior with a simple example using canonical structures. Consider the environment $E = d := 0; c := d$, where there is also a structure with projector **proj**. Suppose further that there is a canonical instance i registered for **proj** and d . Then, the algorithm should succeed finding a solution for the following equation:

$$\mathbf{match} \ c \ \mathbf{with} \ 0 \Rightarrow d \mid _ \Rightarrow 1 \ \mathbf{end} \approx \mathbf{proj} \ ?f \quad (2)$$

where $?f$ is an unknown instance of the structure. More precisely, we expect the left-hand side to be reduced as

$$d \approx \mathbf{proj} \ ?f$$

therefore enabling the use of the canonical instance i to solve $?f$.

This is done in the rule CASE- ιL by weak-head normalizing the left-hand side using the standard $\beta\zeta\delta\Sigma\iota$ reduction rules plus a new reduction rule, θ , which weak-head normalizes scrutinees (Figure 5). Note that we really need this new reduction rule: we cannot consider weak-head reducing the term using δE , as it will destroy the constant d in the example above, nor restrict reduction of the scrutinee to not include δE , as it will be too restrictive (disallowing δE in the reduction on the l.h.s. makes Equation 2 not unifiable).

In Equation 2 we have a **match** on the l.h.s., and a constant on the r.h.s. (the projector). By giving priority to the ι reduction strategy over the δE one we can be sure that the projector will not get unfolded beforehand, and therefore the canonical instance resolution mechanism will work as expected. Different is the situation when we have constants on both sides of the equation. For instance, consider the following equation:

$$c \approx \mathbf{proj} \ ?f \quad (3)$$

in the same context as before. Since there is no instance defined for c , we expect the algorithm to unfold it, uncovering the constant d . Then, it should solve the equation, as before, by instantiating $?f$ with i . If the projector is unfolded first instead, then the algorithm will not find the solution. The reason is that the projector unfolds to a case on the unknown $?f$:

$$c \approx \mathbf{match} \ ?f \ \mathbf{with} \ \text{Constr } a_1 \dots a_n \Rightarrow a_j \ \mathbf{end}$$

(Assuming the projector `proj` corresponds to the j -th field in the structure, and `Constr` is the constructor of the structure.) Now the canonical instance resolution will fail to see that the right-hand side is (was) a projector, so after unfolding `c` and `d` on the left, the algorithm will give up and fail.

In this case we cannot just simply rely on the ordering of rules, since that would make the algorithm sensitive to the position of the terms. In order to solve Equation 3 above, for instance, we need to prioritize reduction on the l.h.s. over the r.h.s., but this prioritization will have a negative impact on equations having the projector on the left instead of the right. The solution is to unfold a constant on the r.h.s. *only if the term does not “get stuck”*, that is, does not evaluate to certain values, like an irreducible **match**. More precisely, we define the concept of “being stuck” as

$$\text{is_stuck } t = \exists t' t''. t \rightsquigarrow_{\delta E, \delta \Gamma}^{0..1} t' \wedge t' \downarrow_{\beta \zeta \iota \theta}^w t'' \text{ and the head of } t'' \text{ is a variable, case, fix, or abstraction}$$

That is, after performing an (optional) δE or $\delta \Gamma$ step and $\beta \zeta \iota \theta$ -weak-head reducing the definition, the head element of the result is tested to be a **match**, **fix**, variable, or a λ -abstraction. Note that the reduction will effectively stop at the first head constant, without unfolding it further. This is important, for instance, when having a definition that reduces to a projector of a structure. If the projector is not exposed, and is instead reduced, then some canonical solution may be lost.

The rule `CONS- δ NOTSTUCKR` unfolds the right-hand side constant only if it will not get stuck. If it is stuck, then the rule `CONS- δ STUCKL` triggers and unfolds the left-hand side, which is precisely what happened in the example above. The rules `CONS- δ` are triggered as a last resort. This controlled unfolding of constants, together with canonical structures resolution, is what allows the encoding of sophisticated meta-programming idioms in [12].

When none of the rules above applies, the algorithm tries η -expansion (`LAM- η` rules). These rules unify a function $\lambda x : T. t$ with a term u . The first premise ensures that u 's head is not an abstraction to avoid overlapping with the `LAM-SAME` rules, otherwise it is possible to build an infinite loop together with the rules `LAM- β` . The following two hypotheses ensure that u has product type with T as domain. First, the type of u is computed as U , and then we ensure U is a product with domain T calling the following function:

$$\begin{aligned} \text{ensure_product}(\bar{\ell} \models C; \Sigma_0; \Gamma; T; U) &= (\phi_2; \Sigma_2) \\ \text{where } \phi_1 &= \bar{\ell}, i \models C \text{ for fresh universe level } i \\ \text{and } \Sigma_1 &= \Sigma_0, ?v : \text{Type}(i)[\Gamma, y : T] \text{ for fresh } ?v \\ \text{and } \phi_1; \Sigma_1; \Gamma \vdash U &\approx_{\equiv} \forall y : T. ?v[\widehat{\Gamma}, y] \triangleright \phi_2; \Sigma_2 \end{aligned}$$

This function returns the result of unifying U with a product type with domain T and unknown range $?v$. For this, the meta-context Σ_0 is extended with $?v$ having type $\text{Type}(i)$, for fresh universe level i , and context Γ extended with $y : T$.

We consider η -expansion as a last resort in the hope of obtaining a function beforehand, after expanding some definition.

4.3 Meta-Variable Instantiation

The rules for meta-variable instantiation are considered in Figure 6, most of which are inspired by Abel & Pientka [1]. There are, however, several differences between their work and ours, since we have a different base logic (CIC instead of LF), and a different assumption on the types of the terms: they require the terms being unified to have the same (known) type, while we do not (types play no directing role in our unification judgment).

For presentation purposes, we only present the rules having a meta-variable on the right-hand side of the equation, but the algorithm also includes the rules with the terms swapped.

Same Meta-Variable: If both terms are the same meta-variable $?x$, we have two distinct cases: if their substitution is exactly the same, the rule `META-SAME-SAME` applies, in which the arguments of the meta-variable are compared point-wise. Note that we require the elements in the substitution to be the same, and not just convertible. If, instead, their substitutions are different, the rule `META-SAME` is attempted. To better understand this rule, let's look at an example. Suppose $?z$ has type $T[x_1 : \text{nat}, x_2 : \text{nat}]$ and we have to solve the equation

$$?z[y_1, y_2] \approx ?z[y_1, y_3]$$

where y_1, y_2 and y_3 are distinct variables. From this equation we cannot know yet what value $?z$ will hold, but at least we know it cannot refer (up to conversion) to the second parameter, x_2 , since then the above equation would have no solution. This reasoning is reflected in the rule `META-SAME` in the hypothesis

$$\Psi_1 \vdash \sigma \cap \sigma' \triangleright \Psi_2$$

This judgment performs an intersection of both substitutions, filtering out those positions from the context of the meta-variable Ψ_1 where the substitutions disagree, resulting in Ψ_2 .

The intersection judgment is defined in the `INTERSEC-*` rules in the same figure. The text inside grey boxes defines a different rule: it allows us to collapse a rule for declarations and a rule for definitions into one only rule. The judgment is conservative: it only filters out different *variables*. The judgment is undefined if the two substitutions have different *terms* (not variables) in some position. Of course, a more aggressive approach is possible, checking for convertibility of the terms instead of just syntactic equality, but it is not clear whether the few more cases covered compensates for the potential loss in performance.

Coming back to the rule `META-SAME`, by filtering out the disagreeing positions of the substitution, we obtain a new context Ψ_2 , which is a subset of Ψ_1 . Since this smaller context may be ill-formed, we *sanitize* it. The sanitization judgment is defined at the bottom of the figure, and it simply removes every (possibly defined) variable in the context whose free variables are not included in the context. This process results in a new (possibly smaller) context Ψ_3 . After making sure that the type T of $?x$ is still well-formed in this context, we restrict $?x$ to only refer to the variables in Ψ_3 . We do this by creating a new meta-variable $?y$ with the type of $?x$, but in the context Ψ_3 . We further instantiate $?x$ with $?y$. Both the creation of $?y$ and the instantiation of $?x$ in the meta-context Σ is expressed in the fragment $\Sigma \cup \{?y : T[\Psi_3], ?x := ?y[\widehat{\Psi}_3]\}$ of the last hypothesis. We use this new meta-context to compare point-wise the arguments of the meta-variable.

Meta-Variable Instantiation: The `META-INST` rules instantiate a meta-variable applying a variation of *higher-order pattern unification* (HOPU) [14]. They unify a meta-variable $?x$ with some term t , obtaining a MGU (actually, as we will see in §4.3.1, *almost* a MGU). As required by HOPU, the meta-variable is applied to a suspended substitution mapping variables to variables, ξ , and a spine of arguments ξ' , of variables only. Assuming $?x$ has (contextual) type $T[\Psi]$, this rule must find a term t''' to instantiate $?x$ such that

$$t \approx ?x[\xi] \xi'$$

that is, after performing the suspended substitution ξ and applying arguments ξ' (formally, $t''' \{ \xi / \widehat{\Psi} \} \xi'$), results in a term convertible to t .

Having contexts Σ_0 and Γ , the new term t''' is crafted from t following these steps:

| | | |
|---|---|--|
| $\frac{\text{META-SAME-SAME} \quad \Sigma; \Gamma \vdash \bar{t} \approx_{\equiv} \bar{u} \triangleright \Sigma'}{\Sigma; \Gamma \vdash ?x[\sigma] \bar{t} \approx_{\mathcal{R}} ?x[\sigma] \bar{u} \triangleright \Sigma'}$ | $\frac{\text{META-SAME} \quad \begin{array}{l} ?x : T[\Psi_1] \in \Sigma \quad \Psi_1 \vdash \sigma \cap \sigma' \triangleright \Psi_2 \quad \cdot \vdash \text{sanitize}(\Psi_2) \triangleright \Psi_3 \\ \text{FV}(T) \subseteq \Psi_3 \quad \Sigma \cup \{?y : T[\Psi_3], ?x := ?y[\widehat{\Psi}_3]\}; \Gamma \vdash \bar{t} \approx_{\equiv} \bar{u} \triangleright \Sigma' \end{array}}{\Sigma; \Gamma \vdash ?x[\sigma] \bar{t} \approx_{\mathcal{R}} ?x[\sigma'] \bar{u} \triangleright \Sigma'}$ | |
| $\frac{\text{META-INSTR} \quad \begin{array}{l} ?x : T[\Psi] \in \Sigma_0 \quad t', \xi_1 = \text{remove_tail}(t; \xi') \quad t' \downarrow_{\beta}^w t'' \quad \Sigma_0 \vdash \text{prune}(?x; \xi, \xi_1; t'') \triangleright \Sigma_1 \quad \Sigma_1; \Gamma \vdash \xi_1 : \bar{U} \\ t''' = \lambda y : U\{\xi, \xi_1/\widehat{\Psi}, \bar{y}\}^{-1}. \Sigma_1(t''')\{\xi, \xi_1/\widehat{\Psi}, \bar{y}\}^{-1} \quad \Sigma_1; \Psi \vdash t''' : T' \quad \Sigma_1; \Psi \vdash T' \approx_{\leq} T \triangleright \Sigma_2 \quad ?x \notin \text{FMV}(t''') \end{array}}{\Sigma_0; \Gamma \vdash t \approx_{\mathcal{R}} ?x[\xi] \xi' \triangleright \Sigma_2 \cup \{?x := t'''\}}$ | | |
| $\frac{\text{META-FOR} \quad \begin{array}{l} ?x : T[\Psi] \in \Sigma_0 \quad 0 < n \quad \Sigma_0; \Gamma \vdash u \overline{u'_m} \approx_{\equiv} ?x[\sigma] \triangleright \Sigma_1 \quad \Sigma_1; \Gamma \vdash \overline{u''_n} \approx_{\equiv} \bar{t}_n \triangleright \Sigma_2 \end{array}}{\Sigma_0; \Gamma \vdash u \overline{u'_m u''_n} \approx_{\mathcal{R}} ?x[\sigma] \bar{t}_n \triangleright \Sigma_2}$ | | |
| $\frac{\text{META-DELDEPSR} \quad \begin{array}{l} ?x : T[\Psi] \in \Sigma \quad l = [i \sigma_i \text{ is variable and } \nexists j > i. \sigma_i = (\sigma, \bar{u})_j] \\ \cdot \vdash \text{sanitize}(\Psi_l) \triangleright \Psi' \quad \text{FV}(T) \subseteq \Psi' \quad \Sigma \cup \{?y : T[\Psi'], ?x := ?y[\widehat{\Psi}']\}; \Gamma \vdash t \approx_{\mathcal{R}} ?y[\sigma_l] \bar{u} \triangleright \Sigma' \end{array}}{\Sigma; \Gamma \vdash t \approx_{\mathcal{R}} ?x[\sigma] \bar{u} \triangleright \Sigma'}$ | | |
| $\frac{\text{META-REDUCER} \quad \begin{array}{l} ?u : T[\Psi] \in \Sigma_0 \quad t \rightsquigarrow_{\delta}^{w, 0.1} t' \quad t' \downarrow_{\beta \zeta_i \theta}^w t'' \quad \Sigma_0; \Gamma \vdash t'' \approx_{\mathcal{R}} ?u[\sigma] \bar{t}_n \triangleright \Sigma_1 \end{array}}{\Sigma_0; \Gamma \vdash t \approx_{\mathcal{R}} ?u[\sigma] \bar{t}_n \triangleright \Sigma_1}$ | $\frac{\text{INTERSEC-NIL}}{\cdot \vdash \cdot \cap \cdot \triangleright \cdot}$ | |
| $\frac{\text{INTERSEC-KEEP} \quad \Psi \vdash \sigma \cap \sigma' \triangleright \Psi'}{\Psi, x := u : A \vdash \sigma, t \cap \sigma', t \triangleright \Psi', x : A}$ | $\frac{\text{INTERSEC-REMOVE} \quad \begin{array}{l} \Psi \vdash \sigma \cap \sigma' \triangleright \Psi' \quad y \neq z \\ \Psi, x := u : T \vdash \sigma, y \cap \sigma', z \triangleright \Psi' \end{array}}{\Psi, x := u : T \vdash \sigma, y \cap \sigma', z \triangleright \Psi'}$ | $\frac{\text{SANITIZE-NIL}}{\xi \vdash \text{sanitize}(\cdot) \triangleright \cdot}$ |
| $\frac{\text{SANITIZE-KEEP} \quad \begin{array}{l} \text{FV}(T) \subseteq \xi \quad \text{FV}(u) \subseteq \xi \quad x, \xi \vdash \text{sanitize}(\Psi) \triangleright \Psi' \end{array}}{\xi \vdash \text{sanitize}(x := u : T, \Psi) \triangleright x : T, \Psi'}$ | $\frac{\text{SANITIZE-REMOVE} \quad \begin{array}{l} \text{FV}(T) \not\subseteq \xi \vee \text{FV}(u) \not\subseteq \xi \quad \xi \vdash \text{sanitize}(\Psi) \triangleright \Psi' \end{array}}{\xi \vdash \text{sanitize}(x := u : T, \Psi) \triangleright \Psi'}$ | |

Figure 6. Meta-variable instantiation.

1. To avoid unnecessarily η -expanded solutions, the term t and arguments ξ' are decomposed using the function `remove_tail(\cdot ; \cdot)`:

$$\begin{aligned} \text{remove_tail}(t \ x; \xi, x) &= \text{remove_tail}(t; \xi) && \text{if } x \notin \text{FV}(t) \\ \text{remove_tail}(t; \xi) &= (t, \xi) && \text{in any other case} \end{aligned}$$

This function, applied to t and ξ' , returns a new term t' and a list of variables ξ_1 , where there exists ξ_2 such that $t = t' \xi_2$ and $\xi' = \xi_1, \xi_2$, and ξ_2 is the longest such list. For instance, in the following example

$$?f[] \ x \ y \approx \text{addn} \ x \ y$$

where `addn` is the addition operation on natural numbers, we want to remove “the tail” on both sides of the equation, leading to the natural solution `?f[] := addn`. In this example, ξ_1 is the empty list, ξ_2 is $[x, y]$, and t' is `addn`.

The check that $x \notin \text{FV}(t)$ in the first case above ensures that no solutions are erroneously discarded. Consider the following equation:

$$?f[] \ x \approx \text{addn0} \ x \ x$$

If we remove the argument of the meta-variable, we will end up with the unsolvable equation `?f[] \approx addn0 x` .

2. The term obtained in the previous step is weak head β normalized, noted $t' \downarrow_{\beta}^w t''$. This is performed in order to remove false dependencies, like variable x in $(\lambda y. 0) \ x$.

3. The meta-variables in t'' are *pruned*. This process is quite involved, and detailed examples can be found in [1]. The formal description will be discussed below in §4.3.1.

At high level, the pruning judgment ensures that the term t'' has no “offending variables”, that is, free variables outside of those occurring in the substitution ξ, ξ_1 . It does so by removing elements from the suspended substitutions occurring in t'' , containing variables outside of ξ, ξ_1 . For instance, in the example `?f[] $x \approx$ addn0 ? $u[x, y]$` , the variable y has to be removed from the substitution on the r.h.s. since it does not occur in the l.h.s. Similarly, if the meta-variable being instantiated occurs inside a suspended substitution, it has to be removed from the substitution to avoid a circularity in the instantiation. The output of this judgment is a new meta-context Σ_1 .

4. The final term t''' is constructed as

$$\lambda y : U\{\xi, \xi_1/\widehat{\Psi}, \bar{y}\}^{-1}. \Sigma_1(t''')\{\xi, \xi_1/\widehat{\Psi}, \bar{y}\}^{-1}$$

First, note that t''' has to be a function taking n arguments \bar{y} , where $n = |\xi_1|$. For the moment, let’s forget about the types of each y_j .

The body of this function is the term obtained from the second step, t'' , after performing a few changes. First, all of its defined meta-variables are normalized with respect to the meta-context obtained in the previous step, Σ_1 , in order to replace the meta-variables with the pruned ones. This step effectively removes false dependencies on variables not occurring in ξ, ξ_1 .

Then, the *inversion* of substitution $\xi, \xi_1 / \hat{\Psi}, \bar{y}$ is performed. This inversion ensures that all free variables in $\Sigma_1(t'')$ are replaced by variables in Ψ and \bar{y} . More precisely, it replaces every variable in $\Sigma_1(t'')$ appearing *only once* in the image of the substitution (ξ, ξ_1) by the corresponding variable in the domain of the substitution $(\hat{\Psi}, \bar{y})$. If a variable appears multiple times in the image and occur in term t'' , then inversion fails.

The type of each argument y_j of the function is the type U_j , obtained from the j -th element in ξ_1 , after performing the inversion substitution (with the caveat that the substitution includes only the $j - 1$ elements in \bar{y}).

5. The type of t''' , which now only depends on the context Ψ , is computed as T' , and unified with the type of $?x$, obtaining a new meta-context Σ_2 .

In the special case where t''' is itself a meta-variable of type an arity (an n -ary dependent product whose codomain is a sort), we do not directly force the type of the instance T' to be smaller than T , which would unnecessarily restrict the universe graph. Instead, we downcast T and T' to a smaller type according to the cumulativity relation before converting them. The idea is that, if we are unifying meta-variables $?x$ and $?y$, with $?x : \text{Type}(i)[\Gamma]$ and $?y : \text{Type}(j)[\Gamma']$, the body of $?x$ and $?y$ just has to be of type $\text{Type}(k)$ for some $k \leq i, j$.

6. Finally, an *occurs check* is performed to prevent illegal solutions, making sure $?x$ does not occur in t''' .

The algorithm outputs Σ_2 plus the instantiation of $?x$ with t''' .

First-Order Approximation: The rules META-INST only applies if the spine of arguments of the meta-variable only have variables. This can be quite restrictive. Consider for instance the following equation that tries to unify an unknown function, applied to an unknown argument, with the term 1 (expanded to S 0):

$$S\ 0 \approx ?f[]\ ?y[]$$

As usual, such equations have multiple solutions, but there is one that is “more natural”: assign S to $?f$ and 0 to $?y$. However, since the argument to the meta-variable is not a variable, it does not comply with HOPU, and therefore is not considered by the META-INST rules. In an scenario like this, the META-FO rules perform a *first order approximation*, unifying the meta-variable ($?f$ in the equation above) with the term on the l.h.s. without the last n arguments (S), which are in turn unified pointwise with the n arguments in the spine of the meta-variable (0 and $?y$, respectively). Note that the rule APP-FO does not subsume this rule, as it requires both terms being equated to have the same number of arguments.

Meta-Variable Dependencies Erasure: If none of the rules above work, the algorithm makes a somewhat brutal attempt: the rule META-DELDepsR chops off every element in the substitution that is not a variable, or that is a duplicated variable. Therefore, problems not complying with HOPU can be reconsidered. One issue with this rule is that it fixes a solution where many solutions may exist. Although the selected solution works most of the time, as we are going to see in §5, it might not be the one expected by the user.

Formally, this rule first takes each position i in σ such that σ_i is a variable with no duplicated occurrence in σ, \bar{u} . The resulting list l containing those positions is used to filter out the local context of the meta-variable, obtaining the new context Ψ' . After sanitizing this context, a fresh meta-variable $?y$ is created in this restricted local context, and $?x$ is instantiated with this meta-variable. The new meta-context obtained after this instantiation is used to recursively call the unification algorithm to solve the problem $?y[\sigma_i]\ \bar{u} \approx t$.

Eliminating Dependencies via Reduction: Sometimes the term being assigned to the meta-variable has variables not occurring

in the substitution, but that can be eliminated via reduction. For instance, take the following equation

$$\pi_1(0, x) \approx ?g[]$$

It has a solution, after reducing the term on the l.h.s., obtaining the easily solvable equation $0 \approx ?g[]$. This is precisely what rules META-REDUCE do, as a last attempt to make progress.

4.3.1 Pruning

Figure 7 shows the actual process of pruning. The pruning judgment is noted

$$\Sigma \vdash \text{prune}(?x; \xi; t) \triangleright \Sigma'$$

It takes a meta-context Σ , a meta-variable $?x$, a list of variables ξ , the term to be pruned t , and returns a new meta-context Σ' , which is an extension of Σ where all the meta-variables with offending variables in their suspended substitution are instantiated with pruned ones.

For brevity, we only show rules for the Calculus of Constructions fragment of CIC, *i.e.*, without considering pattern matching and fixpoints. The missing rules are easy to extrapolate from the given ones. The only interesting case is when the term t is a meta-variable $?y$ applied to the suspended substitution σ . We have two possibilities: either every variable from every term in σ is included in ξ , in which case we do not need to prune (PRUNE-META-NOPRUNE), or there exists some terms which have to be removed (pruned) from σ (PRUNE-META).

These two rules use an auxiliary judgment to prune the local context of the meta-variable Ψ_0 . This judgment has the form

$$\Psi \vdash \text{prune_ctx}(?x; \xi; \sigma) \triangleright \Psi'$$

Basically, it filters out every variable in Ψ where σ has an *offending term*, that is, a term with a free variable not in ξ , or having $?x$ in the set of free meta-variables. Ψ' is the result of this process.

Coming back to the rules in Figure 7, in PRUNE-META-NOPRUNE we have the condition that the pruning of context Ψ_0 resulted in the same context (no need for a change). More interestingly, when the pruning of Ψ_0 results in a new context Ψ_1 , PRUNE-META does the actual pruning of $?y$. Similarly to the rule META-SAME, it first sanitizes the new context Ψ_1 , obtaining a new context Ψ_2 , then it ensures that the type T is valid in Ψ_2 , by pruning variables outside Ψ_2 , and finally instantiates the meta-variable $?y$ with a fresh meta-variable $?z$, having contextual type $T[\Psi_2]$.

It is important to note that, due to conversion, the process of pruning may loose solutions. For instance, consider the following equation:

$$\pi_1(0, ?x[n]) \approx ?y[]$$

The pruning algorithm will remove n from $?x$, although another solution exists by reducing the l.h.s., assigning 0 to $?y$.

4.4 Canonical Structures Resolution

When an instance i of a structure is declared **Canonical**, COQ will add, for each projector, a record in the *canonical structures database* (Δ_{db}). Each record registers a key consisting of the projector p and the head constructor h of the value for that projector in the instance, and a value, the instance i itself. Then, at high level, when the algorithm has to solve an equation of the form $h\ t \approx p\ ?x$, it searches for the key (p, h) in the database, finding that $?x$ should be instantiated with i .

The process is formally described in Figure 8. We always start from an equation of the form:

$$t'' \approx p_j[\bar{\kappa}]\ \bar{p}\ i\ \bar{t}$$

where p_j is a projector of a structure applied to some universe instance $\bar{\kappa}$, \bar{p} are the parameters of the structure, i is the instance

| | | |
|---|---|--|
| $\frac{\text{PRUNE-CONSTANT} \quad h \in \mathcal{S} \cup \mathcal{C}}{\Sigma \vdash \text{prune}(?x; \xi; h) \triangleright \Sigma}$ | $\frac{\text{PRUNE-VAR} \quad x \in \xi}{\Sigma \vdash \text{prune}(?x; \xi; x) \triangleright \Sigma}$ | $\frac{\text{PRUNE-LAM, PRUNE-PROD} \quad \Pi \in \{\lambda, \forall\} \quad \Sigma \vdash \text{prune}(?x; \xi; x; t) \triangleright \Sigma'}{\Sigma \vdash \text{prune}(?x; \xi; \Pi x. t) \triangleright \Sigma'}$ |
| $\frac{\text{PRUNE-LET} \quad \begin{array}{l} \Sigma_0 \vdash \text{prune}(?x; \xi; t_2) \triangleright \Sigma_1 \\ \Sigma_1 \vdash \text{prune}(?x; \xi; x; t_1) \triangleright \Sigma_2 \end{array}}{\Sigma_0 \vdash \text{prune}(?x; \xi; \text{let } x := t_2 \text{ in } t_1) \triangleright \Sigma_2}$ | $\frac{\text{PRUNE-APP} \quad \begin{array}{l} \Sigma_0 \vdash \text{prune}(?x; \xi; t) \triangleright \Sigma_1 \\ \Sigma_i \vdash \text{prune}(?x; \xi; t_i) \triangleright \Sigma_{i+1} \quad i \in [1, n] \end{array}}{\Sigma_0 \vdash \text{prune}(?x; \xi; t \bar{t}_n) \triangleright \Sigma_{n+1}}$ | $\frac{\text{PRUNE-META-NOPRUNE} \quad \begin{array}{l} ?y : T[\Psi_0] \in \Sigma \quad ?x \neq ?y \\ \Psi_0 \vdash \text{prune.ctx}(?x; \xi; \sigma) \triangleright \Psi_0 \end{array}}{\Sigma \vdash \text{prune}(?x; \xi; ?y[\sigma]) \triangleright \Sigma}$ |
| $\frac{\text{PRUNE-META} \quad \begin{array}{l} ?y : T[\Psi_0] \in \Sigma \quad ?x \neq ?y \quad \Psi_0 \vdash \text{prune.ctx}(?x; \xi; \sigma) \triangleright \Psi_1 \\ \cdot \vdash \text{sanitize}(\Psi_1) \triangleright \Psi_2 \quad \Sigma \vdash \text{prune}(?x; \widehat{\Psi}_2; T) \triangleright \Sigma' \end{array}}{\Sigma \vdash \text{prune}(?x; \xi; ?y[\sigma]) \triangleright \Sigma', ?z : T[\Psi_2] \cup \{?y := ?z[\widehat{\Psi}_2]\}}$ | | $\frac{\text{PRUNECTX-NIL}}{\cdot \vdash \text{prune.ctx}(?x; \xi; \cdot) \triangleright \cdot}$ |
| $\frac{\text{PRUNECTX-NOPRUNE} \quad \begin{array}{l} \text{FV}(t) \subseteq \xi \quad ?x \notin \text{FMV}(t) \quad \Psi \vdash \text{prune.ctx}(?x; \xi; \sigma) \triangleright \Psi' \\ \Psi, z : A \vdash \text{prune.ctx}(?x; \xi; \sigma, t) \triangleright \Psi', z : A \end{array}}{\Psi, z : A \vdash \text{prune.ctx}(?x; \xi; \sigma, t) \triangleright \Psi', z : A}$ | $\frac{\text{PRUNECTX-PRUNE} \quad \begin{array}{l} \text{FV}(t) \not\subseteq \xi \vee ?x \in \text{FMV}(t) \quad \Psi \vdash \text{prune.ctx}(?x; \xi; \sigma) \triangleright \Psi' \\ \Psi, x : A \vdash \text{prune.ctx}(?x; \xi; \sigma, t) \triangleright \Psi' \end{array}}{\Psi, x : A \vdash \text{prune.ctx}(?x; \xi; \sigma, t) \triangleright \Psi'}$ | |

Figure 7. Pruning of meta-variables.

(usually a meta-variable), and \bar{t} are the arguments of the projected value, in the case when it has product type. In order to solve this equation the algorithm proceeds as follows:

LOOKUP-CS

$$\frac{\begin{array}{l} (p_j, h, c_\iota) \in \Delta_{\text{db}} \\ \Phi_1, \iota = \text{fresh}(\Phi_0, c_\iota) \quad \iota \rightsquigarrow_{\delta E} \lambda x : \bar{T}. k[\bar{\kappa}'] \bar{p}' \bar{v} \\ \Sigma_1 = \Sigma_0, ?y : \bar{T} \quad \Phi_1 \models \bar{\kappa} = \bar{\kappa}' \triangleright \Phi_2 \\ \Phi_2; \Sigma_1; \Gamma \vdash \bar{p} \approx_{\equiv} p' \{?y/\bar{x}\} \triangleright \Phi_3; \Sigma_2 \end{array}}{\Phi_0; \Sigma_0 \vdash (p_j, \bar{\kappa}, \bar{p}, h) \in ? \Delta_{\text{db}} \triangleright \Phi_3, \Sigma_2, \iota ?y, v_j \{?y/\bar{x}\}}$$

CS-CONSTR

$$\frac{\begin{array}{l} \Phi_0; \Sigma_0 \vdash (p_j, \bar{\kappa}, \bar{p}, c) \in ? \Delta_{\text{db}} \triangleright \Phi_1, \Sigma_1, \iota, c[\bar{\ell}'] \bar{u}' \\ \Phi_1 \models \bar{\ell} = \bar{\ell}' \triangleright \Phi_2 \quad \Phi_2; \Sigma_1; \Gamma \vdash \bar{u} \approx_{\equiv} \bar{u}' \triangleright \Phi_3; \Sigma_2 \\ \Phi_3; \Sigma_2; \Gamma \vdash i \approx_{\equiv} \iota \triangleright \Phi_4; \Sigma_3 \\ \Phi_4; \Sigma_4; \Gamma \vdash \bar{t}' \approx_{\equiv} \bar{t} \triangleright \Phi_5; \Sigma_4 \end{array}}{\Phi_0; \Sigma_0; \Gamma \vdash c[\bar{\ell}] \bar{u} \bar{t}' \approx_{\mathcal{R}} p_j[\bar{\kappa}] \bar{p} i \bar{t} \triangleright \Phi_5; \Sigma_4}$$

CS-PRODR

$$\frac{\begin{array}{l} \Phi_0; \Sigma_0 \vdash (p_j, \bar{\kappa}, \bar{p}, \rightarrow) \in ? \Delta_{\text{db}} \triangleright \Phi_1, \Sigma_1, \iota, u \rightarrow u' \\ \Phi_1; \Sigma_1; \Gamma \vdash t \approx_{\equiv} u \triangleright \Phi_2; \Sigma_2 \\ \Phi_2; \Sigma_2; \Gamma \vdash t' \approx_{\mathcal{R}} u' \triangleright \Phi_3; \Sigma_3 \\ \Phi_3; \Sigma_3; \Gamma \vdash i \approx_{\equiv} \iota \triangleright \Phi_4; \Sigma_4 \end{array}}{\Phi_0, \Sigma_0; \Gamma \vdash t \rightarrow t' \approx_{\mathcal{R}} p_j[\bar{\kappa}] \bar{p} i \triangleright \Phi_4; \Sigma_4}$$

CS-SORTR

$$\frac{\begin{array}{l} \Phi_0; \Sigma_0 \vdash (p_j, \bar{\kappa}, \bar{p}, s) \in ? \Delta_{\text{db}} \triangleright \Phi_1, \Sigma_1, \iota, v_j \\ \Phi_1; \Sigma_1; \Gamma \vdash s \approx_{\mathcal{R}} v_j \triangleright \Phi_2; \Sigma_2 \\ \Phi_2; \Sigma_2; \Gamma \vdash i \approx_{\equiv} \iota \triangleright \Phi_3; \Sigma_3 \end{array}}{\Phi_0; \Sigma_0; \Gamma \vdash s \approx_{\mathcal{R}} p_j[\bar{\kappa}] \bar{p} i \triangleright \Phi_3; \Sigma_3}$$

CS-DEFAULTR

$$\frac{\begin{array}{l} \Phi_0; \Sigma_0 \vdash (p_j, \bar{\kappa}, \bar{p}, -) \in ? \Delta_{\text{db}} \triangleright \Phi_1, \Sigma_1, \iota, v_j \\ \Phi_3; \Sigma_2; \Gamma \vdash t \approx_{\mathcal{R}} v_j \triangleright \Phi_4; \Sigma_3 \\ \Phi_4; \Sigma_3; \Gamma \vdash i \approx_{\equiv} \iota \triangleright \Phi_5; \Sigma_4 \end{array}}{\Phi_0; \Sigma_0; \Gamma \vdash t \approx_{\mathcal{R}} p_j[\bar{\kappa}] \bar{p} i \triangleright \Phi_5; \Sigma_4}$$

Figure 8. Canonical structures resolution.

1. First, a constant c_ι is selected from Δ_{db} , keying on the projector p_j and the head element h of t'' . The constant c_ι stored in the database is a potentially polymorphic constant, so we build a fresh instance of it (ι) and add the fresh universe levels and constraints to the universe context. Its body is a function taking arguments $\bar{x} : \bar{T}$ and returning the term $k[\bar{\kappa}'] \bar{p}' \bar{v}$, with k the constructor of the structure applied to a universe instance $\bar{\kappa}'$, \bar{p}' the parameters of the structure, and \bar{v} the values for each of the fields of the structure.
2. Then, the expected and inferred universe instances and parameters of the instance are unified, after replacing every argument x with a corresponding fresh meta-variable $?y$.
3. According to the class of h , the algorithm considers different rules:
 - (a) CS-CONST if h is a constant c .
 - (b) CS-PROD if h is a non-dependent product $t \rightarrow t'$.
 - (c) CS-SORT if h is a sort s .

If these do not apply, then it tries CS-DEFAULT.

4. Next, the term t'' is unified with the corresponding projected term in the value of the instance for the j -th field. If t'' is a constant c applied to arguments \bar{u} and the value v_j of the j -th field of ι is c applied to \bar{u}' , then the universe instances and arguments \bar{u} and \bar{u}' are unified. If t'' is a product with premise t and conclusion t' , they are unified with the corresponding terms (u and u') in v_j . Note that cumulativity is preserved in the codomain of products and in the CS-SORT rule.
5. The instance of the structure i is unified with the instance found in the database, ι , applied to the meta-variables \bar{y} . Typically, i is a meta-variable, and this step results in instantiating the meta-variable with the constructed instance.
6. Finally, for CS-CONST only, if the j -th field of the structure has product type, and is applied to \bar{t}' arguments, then these arguments are unified with the arguments \bar{t} of the projector.

As with the rules for meta-variable instantiation, we only show the rules in one direction, with the projector on the right-hand side, but the algorithm also includes the rules in the opposite direction.

4.5 Rule Priorities and Backtracking

The figures shown above does not precisely nail the priority of the rules, nor when the algorithm backtracks. Below we show the precise order of application of the rules, where the rules in the same line are tried in the given order *without* backtracking (the first one matching the conclusion and whose side-conditions are satisfied is used). Rules in different lines or in the same line separated by | are tried *with* backtracking (if one fails to apply, the next is tried). Note that if at any point the environment and the two terms to be unified are ground (they do not contain meta-variables), unification is skipped entirely and a call to COQ’s efficient conversion algorithm is made instead.

1. If a term has a *defined* meta-variable in its head position:
 - (a) META- δ R, META- δ L
2. If both terms heads are *the same undefined* meta-variable:
 - (a) META-SAME-SAME, META-SAME
3. If one term has an *undefined* meta-variable, and the other term does not have the same meta-variable in its head position:

META-INSTR | META-FOR | META-REDUCER | META-DELDEPSR |
 LAM- η R | META-INSTL | META-FOL | META-REDUCEL |
 META-DELDEPSL | LAM- η L
4. Else:
 - (a) If the two terms have different head constants:
 - i. (CS-CONSTR, CS-PRODR, CS-SORTR) | CS-DEFAULTR
 - ii. (CS-CONSTL, CS-PRODL, CS-SORTL) | CS-DEFAULTL
 - (b) APP-FO
 - (c) The remaining rules from Figure 4 in the following order, backtracking only if the hypotheses that are not recursive calls to the algorithm fail to apply:

LAM- β R | LET- ζ R | CASE- ι R | LAM- β L | LET- ζ L | CASE- ι L |
 CONS- δ NOTSTUCKR | CONS- δ STUCKL | CONS- δ R | CONS- δ L |
 LAM- η R | LAM- η L

4.6 A Deliberate Omission: Constraint Postponement

The technique of *constraint postponement* [18] is widely adopted in unification algorithms, including the current algorithm of COQ. It has however some negative impact in COQ, and, as it turns out, it is not as crucial as generally believed.

First, let us show why this technique is incorporated into proof assistants. Sometimes the unification algorithm is faced with an equation that has multiple solutions, in a context where there should only be one possible candidate. For instance, consider the following term witnessing an existential quantification:

$$\text{exist } _ 0 (\text{le_n } 0) : \exists x. x \leq x$$

where exist is the constructor of the type $\exists x. P x$, with P a predicate over the (implicit) type of x . More precisely, exist takes a predicate P , an element x , and a proof that P holds for x , that is, $P x$. In the example above we are providing an underscore in place of P , since we want COQ to find out the predicate, and we annotate the term with a typing constraint (after the colon) to specify that the whole term is a proof of existence of a number lesser or equal to itself. In this case, we provide 0 as such number, and the proof le_n 0, which has type $0 \leq 0$.

During typechecking, COQ first infers the type of the term on the left of the colon, and only then it verifies that this type is compatible (*i.e.*, unifiable) with the typing constraint. When inferring the type for the term on the left, COQ will create a fresh meta-variable for the predicate P , let’s call it $?P$, and unify $?P 0$ with $0 \leq 0$, the type of le_n 0. Without any further information, COQ has four different (incomparable) solutions for P : $\lambda x. 0 \leq 0$, $\lambda x. x \leq 0$, $\lambda x. 0 \leq x$, $\lambda x. x \leq x$.

When faced with such an ambiguity, COQ postpones the equation in the hope that further information will help disambiguate the problem. In this case, the necessary information is given later on through the typing constraint, which narrows the set of solutions to a unique solution.

Constraint postponement has its consequences, though: On one hand, the algorithm can solve more unification problems and hence fewer typing annotations are required (*e.g.*, we do not need to specify P). On the other hand, since constraints are delayed, the algorithm becomes hard to debug and, at times, slow. The reason for these assertions comes from the realisation that the algorithm will continue to (try to) unify the terms, piling up constraints on the way, perhaps to later on find out that, after all, the terms are not unifiable (or are unifiable only if some decision is taken on the delayed equations).

When combined with canonical structures resolution, or any other form of proof automation, this technique is particularly bad, as it may break the assumption that certain value has been previously assigned. The motivation to omit this technique came from experience in projects on proof automation by the first author [12, 25], and on bi-directional elaboration by the second author (in the above example, a bi-directional elaboration algorithm will unify the type returned by exist with the expected type, and only then unify the type of its arguments, thereby posing the unification problems in the right order).

Our results (§5) show that this technique is not crucial.

4.7 Correctness

Our algorithm should satisfy the following correctness criterion: if two well-typed terms t_1 and t_2 unify under universe context Φ and meta-context Σ , resulting in a new universe Φ' and meta-context Σ' , both terms should also be well typed under Σ' . Moreover, both terms should be convertible (or in the cumulativity relation) under Φ', Σ' .

However, this is false—for both the current algorithm implemented in COQ, and the one described here. The culprit is the syntactic check required at typechecking to ensure termination of fix-points, the guard condition. Indeed, it is easy to make unification instantiate a meta-variable with a term containing a non-structurally-recursive call to a recursive function in its context, resulting in an ill-typed term. Hence, we must weaken this conjecture to use a weaker notion of typing, as in Coen’s thesis [19].

For the moment we lack a correctness proof. This work sets the first stone presenting a specification faithful to an implementation that performs well on a variety of large examples (§5). We anticipate that the proof will be simpler than for existing algorithms, notably due to the lack of postponement which usually complicates the argument of type preservation.

5. Evaluation of the Algorithm

Since, as we saw in §4.6, our algorithm does not incorporate certain heuristics, it is reasonable to expect that it will fail to solve several unification problems appearing in existing libraries. To test our algorithm “in the wild” we developed a plugin called UniCoq³, which, when requested, changes the current unification algorithm

³Sources can be downloaded from <http://github.com/unicog>.

of COQ with ours. With this plugin, we compiled four different libraries, and evaluated the number of lines that required changes. These changes may be necessary either because UniCoq found a different solution from the expected one, or because it found no solution at all. As it turns out, UniCoq solved most of the problems it encountered.

The first set of files we considered is the standard library of COQ. With UniCoq, it compiles almost out of the box, with only a few lines requiring extra typing annotations. We believe the reason for such success is that most of the files in the library are several years old, and were conceived in older versions of COQ, when it had a much simpler unification algorithm.

The second set of files come from Adam Chlipala’s book “Certified Programming with Dependent Types” (CPDT) [7]. This book provides several examples of functional programming with dependent types, including several non-trivial unification patterns coming from dependent matches. As a result, from a total of 6,200 lines, only 14 required extra typing annotations. It is interesting to note that 8 of those lines are solved with the use of a bi-directional elaboration algorithm [e.g., 4] enabled by COQ’s **Program** keyword. For instance, some lines construct witnesses for existential quantification, similar to the example shown in §4.6.

The third one is the Mathematical Components library [11], version 1.5beta1. This library presents several challenges, making it appealing for our purpose: (1) It is a huge development, with a total of 78 theory files. (2) It uses canonical structures heavily, providing us with several examples of canonical structures idioms that UniCoq should support. (3) It uses its own set of tactics uniformly calling the same unification algorithm used for elaboration. This last point is extremely important, although a bit technical. Truth be told, COQ has actually two different unification algorithms. One of these algorithms is mainly used by elaboration, and it outputs a sound substitution (up to bugs). This is the one mentioned in this paper as “the original unification algorithm of COQ”. The other algorithm is used by most of COQ’s original tactics (like `apply` or `rewrite`), but it is unsound (in COQ 8.4, it may return ill-typed solutions). `Ssreflect`’s tactics use the former algorithm which is the one being replaced by our plugin. From the 82,000 lines in the library, only 40 lines required changes.

The last set of files also focuses in different canonical structures idioms: the files from Lemma Overloading [12]. It compiles almost as-is, with only one line requiring an extra annotation.

5.1 A Word on the META-DELDEPS Rules

In a sense, the rules `META-DELDEPS` are a bit brutal: they fix an arbitrary solution from the set of possible ones, which might not be the one expected. However, as the numbers above suggest, it works most of the time. In this section we analyse, for the Mathematical Components library, the origins of the unification problems that fail when this rule is turned off (totaling +300 lines).

Non-dependent if–then–elses: Most notably, the culprit for about two thirds of the failures are `Ssreflect`’s **if–then–elses**. In `Ssreflect`, the type of the branches of an **if** are assumed to depend on the conditional. For instance, the example **if** b **then** 0 **else** 1 fails to compile if the `Ssreflect` library is imported. With `Ssreflect`, a fresh meta-variable $?T$ is created for the type of the branches, with contextual type `Type[b : true]`. When unifying it with the actual type of each branch, b is substituted by the corresponding boolean constructor. This results in the following equations:

$$?T[\text{true}] \approx \text{nat} \quad ?T[\text{false}] \approx \text{nat}$$

Since they are not of the form required by `HOPU`, UniCoq (without the `META-DELDEPS` rules) fails.

False dependency in the in modifier: Another less common issue comes from the `in` modifier in `Ssreflect`’s `rewrite` tactic. This modifier allows the selection of a portion of the goal to perform the rewrite. For instance, if the goal is $1 + x = x + 1$ and we want to apply commutativity of addition on the term on the right, we can perform the following rewrite:

$$\text{rewrite [in } X \text{ in } _ = X] \text{addnC}$$

With the rule, UniCoq instantiates X with the r.h.s. of the equation, and `rewrite` applies commutativity only to that portion of the goal. Without it, however, `rewrite` fails. In this case, the hole $(_)$ is replaced by a meta-variable $?y$, which is assumed to depend on X . But X is also replaced by a meta-variable, $?z$, therefore the unification problem becomes

$$?y[x, ?z[x]] = ?z[x] \approx 1 + x = x + 1$$

that, in turn, poses the equation $?y[x, ?z[x]] \approx 1 + x$, which does not have an MGU.

Non-dependent products: About 30 lines required a simple typing annotation to remove dependencies in products. Consider the following COQ term:

$$\forall P x. (P (S x) = \text{True})$$

When COQ elaborates this term, it first assigns P and x two unknown types, $?T$ and $?U$ respectively, the latter depending on P . Then, it elaborates the term on the left of the equal sign, obtaining further information about the type $?T$ of P : it has to be a (possibly dependent) function $\forall y : \text{nat}. ?T'[y]$. The type of the term on the left is the type of P applied to $S x$, that is, $?T'[S x]$. After elaborating the term on the right and finding out it is a `Prop`, it unifies the types of the two terms, obtaining the equation

$$?T'[S x] \approx \text{Prop}$$

Since, again, this equation does not comply with `HOPU`, UniCoq fails without `META-DELDEPS`.

Explicit duplicated dependencies: There are 15 occurrences where the proof developer wrote explicitly a dependency that duplicates an existing one. Consider for instance the following rewrite statement:

$$\text{rewrite } [_ + _ w] \text{addnC}$$

Here, the proof developer intends to rewrite using commutativity on a fragment of the goal matching the pattern $_ + _ w$. Let’s assume that in the goal there is one occurrence of addition having w occurring in the right, say $t + (w + u)$, for some terms t and u . Since the holes $(_)$ are elaborated as a meta-variable depending on the entire local context, in this case it will include w . Therefore, the pattern will be elaborated as $?y[w] + ?z[w] w$ (assuming no other variables appear in the local context). When unifying the pattern with the desired occurrence we obtain the problem:

$$?z[w] w \approx w + u$$

This equation does not have a MGU, since either w on the l.h.s. can be used as a representative for the w on the r.h.s.. The rules `META-DELDEPS` remove the inner w .

Looking closely into these issues, it seems as if the dependencies were incorrectly introduced in the first place. We plan to study modifications to elaboration and tactics to avoid these dependencies, and study the impact of such changes.

6. Closing remarks

We presented the first formalization of a realistic unification algorithm for COQ, featuring overloading and universe polymorphism.

Moreover, we give a precise characterization of *controlled backtracking* (rules APP-FO plus CONS-*), which, together with the rules for overloading (Figure 8), allow us to explain the patterns introduced in [12]. The algorithm presented in this work is predictable, in the sense that the order in which subproblems are evaluated can be deduced directly from the rules. In particular, we have not introduced the technique of constraint postponement, which reorders unification subproblems. This omission, made in favor of predictability, has shown not to be problematic in practice (§5).

The algorithm includes a heuristic, incarnated in the rules META-DELDEPS that forces a non-dependent solution where multiple solutions might exist. We have studied various scenarios where it is being used, and shown that this heuristic can be replaced in most cases by smarter tactics and elaboration algorithms (§5.1).

In the future we plan to prove soundness of the algorithm (see §4.7), and to improve its performance to make it significantly faster than the current algorithm of COQ.

Acknowledgments

We are deeply grateful to Georges Gonthier for his suggestion on adding the META-DELDEPS rules, Enrico Tassi for carefully explaining the θ reduction strategy and its use, and Andreas Abel, Derek Dreyer, Hugo Herbelin, Aleksandar Nanevski, Scott Kilpatrick, Viktor Vafeiadis, and the anonymous reviewers for their important feedback on earlier versions of this work. This research was partially supported by EU 7FP grant agreement 295261 (MEALS).

References

- [1] A. Abel and B. Pientka. Higher-order dynamic pattern unification for dependent types and records. In *TLCA*. Springer, 2011.
- [2] A. Asperti, C. S. Coen, E. Tassi, and S. Zacchiroli. Crafting a proof assistant. In *TYPES*. Springer-Verlag, 2006.
- [3] A. Asperti, W. Ricciotti, C. S. Coen, and E. Tassi. Hints in unification. In *TPHOLS*, volume 5674 of *LNCS*. Springer, 2009.
- [4] A. Asperti, W. Ricciotti, C. S. Coen, and E. Tassi. A Bi-Directional Refinement Algorithm for the Calculus of (Co)Inductive Constructions. *LMCS*, 8(1), 2012.
- [5] E. Brady. Idris, a general-purpose dependently typed programming language: Design and implementation. *JFP*, 23, 2013.
- [6] I. Cervesato and F. Pfenning. A linear spine calculus. *Journal of Logic and Computation*, 13(5):639–688, 2003.
- [7] A. Chlipala. *Certified Programming with Dependent Types*. MIT Press, 2011. <http://adam.chlipala.net/cpdt/>.
- [8] F. Garillot. *Generic Proof Tools and Finite Group Theory*. PhD thesis, Ecole Polytechnique X, Dec. 2011.
- [9] F. Garillot, G. Gonthier, A. Mahboubi, and L. Rideau. Packaging Mathematical Structures. In *TPHOL*. Springer, 2009.
- [10] G. Gonthier, A. Asperti, J. Avigad, Y. Bertot, C. Cohen, F. Garillot, S. Le Roux, A. Mahboubi, R. O’Connor, S. Ould Biha, I. Pasca, L. Rideau, A. Solovyev, E. Tassi, and L. Théry. A machine-checked proof of the odd order theorem. In *ITP*. Springer, 2013.
- [11] G. Gonthier, A. Mahboubi, and E. Tassi. A small scale reflection extension for the Coq system. Technical report, INRIA, 2008.
- [12] G. Gonthier, B. Ziliani, A. Nanevski, and D. Dreyer. How to make ad hoc proof automation less ad hoc. *JFP*, 23(04):357–401, 2013.
- [13] A. Mahboubi and E. Tassi. Canonical Structures for the working Coq user. In *ITP*. Springer, 2013.
- [14] D. Miller. Unification of simply typed lamda-terms as logic programming. In *ICLP*. MIT Press, 1991.
- [15] A. Nanevski, F. Pfenning, and B. Pientka. Contextual modal type theory. *ACM Trans. Comput. Logic*, 9(3), June 2008.
- [16] U. Norell. Dependently Typed Programming in Agda. In *TLDI*. ACM, 2009.
- [17] S. Peyton Jones, D. Vytiniotis, S. Weirich, and G. Washburn. Simple unification-based type inference for gadts. In *ICFP*. ACM, 2006.
- [18] J. Reed. Higher-order constraint simplification in dependent type theory. In *LFMTP*, 2009.
- [19] C. Sacerdoti Coen. *Mathematical Knowledge Management and Interactive Theorem Proving*. PhD thesis, University of Bologna, 2004.
- [20] A. Saïbi. *Outils Generiques de Modelisation et de Demonstration pour la Formalisation des Mathematiques en Theorie des Types. Application a la Theorie des Categories*. PhD thesis, University Paris 6, 1999.
- [21] M. Sozeau and N. Tabareau. Universe Polymorphism in Coq. In *ITP*. Springer, 2014.
- [22] The Coq Development Team. *The Coq Proof Assistant Reference Manual – Version V8.4*, 2012.
- [23] P. Wadler and S. Blott. How to make ad-hoc polymorphism less ad hoc. In *POPL*, pages 60–76, 1989.
- [24] B. Ziliani, D. Dreyer, N. Krishnaswami, A. Nanevski, and V. Vafeiadis. Mtac: A monad for typed tactic programming in coq. *To appear in JFP*, ??(?):??–??, 2015.
- [25] B. Ziliani, D. Dreyer, N. R. Krishnaswami, A. Nanevski, and V. Vafeiadis. Mtac: A monad for typed tactic programming in Coq. In *ICFP*, 2013.