

# NORA/HAMMR: Making Deduction-Based Software Component Retrieval Practical\*

Johann Schumann

Automated Reasoning, Institut für Informatik, TU München, D-80290 München

`schumann@informatik.tu-muenchen.de`

Bernd Fischer

Abt. Softwaretechnologie, TU Braunschweig, D-38092 Braunschweig

`fisch@ips.cs.tu-bs.de`

## Abstract

*Deduction-based software component retrieval uses pre- and postconditions as indexes and search keys and an automated theorem prover (ATP) to check whether a component matches. This idea is very simple but the vast number of arising proof tasks makes a practical implementation very hard. We thus pass the components through a chain of filters of increasing deductive power. In this chain, rejection filters based on signature matching and model checking techniques are used to rule out non-matches as early as possible and to prevent the subsequent ATP from “drowning.” Hence, intermediate results of reasonable precision are available at (almost) any time of the retrieval process. The final ATP step then works as a confirmation filter to lift the precision of the answer set. We implemented a chain which runs fully automatically and uses MACE for model checking and the automated prover SETHEO as confirmation filter. We evaluated the system over a medium-sized collection of components. The results encourage our approach.*

## 1. Introduction

Reuse of approved software components has been identified as one of the key factors for successful software engineering projects. Although the reuse process also covers many non-technical aspects [33], retrieving appropriate software components from a reuse library is a central task. This is best captured by the *First Golden Rule for Software*

*Reuse*: “You must find it before you can reuse it!”<sup>1</sup>

Most earlier software component retrieval (SCR) methods (e.g., [19]) grew out of classical information retrieval for unstructured texts. However, since software components are highly structured, more specialized approaches may lead to better results. In this paper we will concentrate on a deduction-based approach where we use pre- and postconditions as the components’ indexes and as search keys. A component matches a search key if the involved pre- and postconditions satisfy a well-defined logical relation, e.g., if the component has a weaker precondition and a stronger postcondition than the search key. From this matching relation a proof task is constructed and an ATP is used to establish (or disprove) the match.

This approach has been proposed before (e.g., [28, 20]) but without convincing success because essential user requirements have been neglected. In this paper we follow a more user-oriented approach and describe steps for making deduction-based SCR practical. We concentrate on deduction-based SCR because it is the key technique which underlies more ambitious logic-based software engineering approaches, e.g., program synthesis [17] or component adaptation [25]. For a discussion of benefits and the integration into software engineering processes we refer to [9].

In the next two sections we outline the user requirements for a practical reuse tool and present our system architecture, featuring the filter pipeline and a graphical user interface. Then, we discuss the construction of proof tasks out of the given VDM-SL specifications. This is an important step as the different approaches model quite different reuse aspects. Sections 5 and 6 focus on the two major components of the filter pipeline, namely rejection of non-matches with model-checking techniques and the final confirmation filter with SETHEO. We have evaluated our approach over

---

\*This work is supported by the DFG within the Schwerpunkt “Deduktion” (grant Sn11/2-3), the habilitation grant Schu908-1/5, and the Sonderforschungsbereich SFB 342, Subproject A5 (PARIS: Parallelization of Inference Systems). Part of the work was done while visiting the ICSI Berkeley.

---

<sup>1</sup>This rule has been attributed to Will Tracz.

a database of 55 list specifications. We present and assess the results of these first experiments in Section 7. Finally, we compare our approach to the related work and conclude with current and future work on NORA/HAMMR<sup>2</sup>.

## 2. The User's Point of View

Most earlier work focussed on the technical aspects of deduction-based SCR. The users had to write complete specifications in the ATP input language and even had to supply useful lemmata. The provers were run in batch mode, checking the whole library before any results were presented. ATP runtimes and problems of scaling-up were ignored.

This view led to severe acceptance problems as the users are software engineers and no ATP experts.<sup>3</sup> Their main requirements are that the tool is easy to use, fully automatized, fast and customizable, and hides all evidence of automated theorem proving. Hiding the ATP has some consequences. The input language must be a fully-flavored specification language and not pure first order logic (FOL). But then the automatic construction of the actual proof tasks becomes itself a major task.

Short response times are also essential as the *Fourth Reuse Truism* demands that “you must find it faster than you can rebuild it!” [15]. However, due to the computational complexity of ATP, truly interactive (“sub-second”) behavior is still far out of reach. Instead, *anytime behavior* is acceptable: intermediate results of sufficient precision must be available to the user at (almost) any time during the retrieval process. Retrieval may then be guided further with feedback from the user who may for example strengthen the search key incrementally.

Ideally, the tool doesn't constrain user feedback to the queries but allows a customization of the complete retrieval process. This includes the selection of an appropriate match relation from a given list of choices as well as some tuning of the deductive mechanism (e.g., time limits or model sizes). But it is important to ensure that the tool still runs fully automatically and produces useful results even without customization.

In exchange for these constraints, deduction-based SCR offers the unique feature that completeness and even soundness are not absolutely vital. Incomplete and unsound deduction methods only reduce *recall* (“do we get all matching components?”) and *precision* (“do we get the right components?”).

<sup>2</sup>NORA is no real acronym, HAMMR is a highly adaptive multi-method retrieval tool.

<sup>3</sup>In a real-life setting, a reuse administrator is required who knows the applied deduction methods and who can “tune” libraries (e.g., by giving ATP settings and domain specific lemmata).

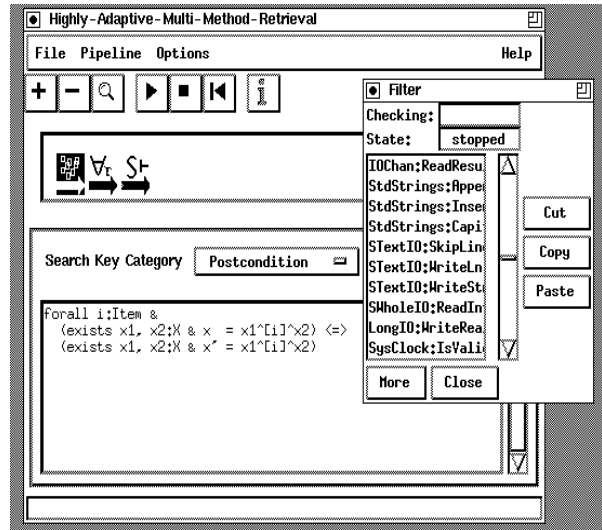


Figure 1. Graphical user interface

## 3. System Architecture

In order to meet the user requirements we implemented NORA/HAMMR as a *filter pipeline* through which the candidates are fed. This pipeline typically starts with *signature matching filters*. They check whether candidate and query have “compatible” calling conventions (i.e., types or signatures). The notion of compatibility is specified by an equational theory  $E$ ; the filter then applies  $E$ -matching or  $E$ -unification of the signature terms. Typical theories include axioms to handle associativity and commutativity of parameter lists and records, currying (for functional languages), pointer types and VAR-parameters (for imperative languages), and coercion rules (see [8] for a detailed discussion).

Then, *rejection filters* try to eliminate non-matches as fast as possible. This is a crucial step to prevent the ATP from “drowning” as there are many more non-matching than matching candidates. We currently apply model generation techniques to check the validity of the tasks in suitable finite models. However, both precision and recall may decrease because this approach is neither sound nor complete.

Finally, *confirmation filters* check the validity of the remaining proof tasks and thus lift the precision of the result to 100%. Here, we apply SETHEO, a high-performance prover based on the Model Elimination calculus. Both filter classes will be described below in more detail.

The graphical user interface (cf. Figure 1) reflects the idea of successive filtering. The pipeline may easily be customized through an icon pad; each filter icon also hides a specialized filter control window which allows some fine-tuning of the filters. Additional inspectors display intermediate results and grant easy access to the components. They

also allow to save intermediate results in a file such that they may easily be used as libraries for subsequent retrieval runs. The objective of the GUI is precisely to hide all evidence of ATP usage. Hence, the knowledge necessary to use NORA/HAMMR as a *tool* is restricted to VDM-SL [7] which we use as our input language and to the target language which is required for signature matching.

#### 4. Proof Tasks and Reuse

The overall structure of the generated proof tasks depends on the definition of the *match relation* which is used in a deduction-based SCR tool. Thus it ultimately depends on the kind of reuse which the tool aims at.

Most often, deduction-based SCR is configured to ensure *plug-in compatibility* of the retrieved components:  $c$  matches if it has a weaker precondition and a stronger postcondition than the search key  $q$ . This is usually (cf. e.g., [34]) formalized as  $(pre_q \Rightarrow pre_c) \wedge (post_c \Rightarrow post_q)$ <sup>4</sup>. However, this is not adequate for partial functions. If  $q$  is a partial function (e.g. *tail*) and  $c$  its total completion (e.g.,  $c(nil)$  returns *nil*) then we want  $c$  to match  $q$  even if its “completed” result does not fit the original specification. It is thus necessary to restrict the implication between the postconditions on the domain given by  $pre_q$ . We thus work with proof tasks of the form

$$(pre_q \Rightarrow pre_c) \wedge (pre_q \wedge post_c \Rightarrow post_q) \quad (1)$$

which are similar to [34]’s “guarded plug-in match” except for our use of the stronger (via the first implication) precondition from the query. Plug-in compatibility supports safe reuse. The retrieved components may be considered as black boxes and may be reused “as is”, without further proviso or modification of the component.

Sometimes plug-in compatibility is not applicable because the users don’t want to specify any precondition but are willing to accept whatever comes, as long as their postconditions are met. In that case, (1) simplifies to

$$pre_c \wedge post_c \Rightarrow post_q \quad (2)$$

or *conditional compatibility*. However, reuse now becomes potentially unsafe because any client still has to satisfy the open obligation  $pre_q$ .

Sometimes (2) might be too strong, and retrieves no components, although the library contains “almost” matches, e.g., partial functions. To additionally retrieve such components, *partial compatibility* may be used:

$$pre_c \wedge pre_q \wedge post_c \Rightarrow post_q \quad (3)$$

<sup>4</sup>Actually, the proof tasks are universally closed wrt. the formal input and output parameters of the component and the query and also contain equations relating the parameters. Likewise, the pre- and postconditions are of course logical functions of the respective parameters. However, to improve readability, we use these traditionally abbreviated formulations.

Thus, a component is retrieved if it computes the correct results on the common domain. If, however, the domains of  $c$  and  $q$  are disjoint,  $pre_c$  and  $pre_q$  are never true at the same time and thus (3) will become vacuously true. But usually  $q$  and  $c$  then also work on different types and  $c$  should already be rejected by signature matching. If  $c$  has an empty domain or is not implementable (i.e.,  $post_c$  never becomes true), (3) will again become vacuously true and  $c$  will be retrieved for any query. However, this should not happen in a well-designed library.

Obviously, reuse based on partial compatibility is unsafe because the retrieved components are not guaranteed to work on the entire required domain. But they might be good starting points for desired more general implementations. Hence, the components must be considered as “white boxes”—their code needs a closer inspection.

As an example, let us consider the following VDM-SL specifications:<sup>5</sup>

$$\begin{array}{ll} rotate(l : List) l' : List & shuffle(x : X) x' : X \\ \text{pre true} & \text{pre true} \\ \text{post } (l = [] \Rightarrow l' = []) \wedge & \text{post } \forall i : Item. \\ (l \neq [] \Rightarrow & (\exists x_1, x_2 : X \cdot x = x_1 \wedge [i] \wedge x_2 \Leftrightarrow \\ l' = (tl\ l) \wedge [hd\ l]) & \exists x_1, x_2 : X \cdot x' = x_1 \wedge [i] \wedge x_2) \end{array}$$

Let us further assume that we use plug-in compatibility as match relation, *rotate* as candidate  $c$  and *shuffle* as query  $q$ . Then several steps are necessary to construct a sorted FOL proof task. First, the formal parameters must be identified, in this case  $l = x$  and  $l' = x'$ <sup>6</sup>. Then, VDM’s underlying three-valued logic LPF must be translated into FOL. This essentially requires the explicit insertion of additional preconditions into the proof task to prevent reasoning from undefined terms as well as a translation of the connectives which takes care of the missing law of the excluded middle [12, 22]. In our example, this results in the proof task

$$\begin{array}{l} \forall l, l', x, x' : List \cdot (l = x \wedge l' = x' \wedge \text{true} \Rightarrow \text{true}) \\ \wedge (l = x \wedge l' = x' \wedge (l = [] \Rightarrow l' = [])) \\ \wedge (l \neq [] \Rightarrow (l \neq [] \Rightarrow l' = (tl\ l) \wedge [hd\ l])) \\ \Rightarrow (\forall i : Item \cdot (\exists x_1, x_2 : X \cdot x = x_1 \wedge [i] \wedge x_2 \\ \Leftrightarrow \exists x_1, x_2 : X \cdot x' = x_1 \wedge [i] \wedge x_2)) \end{array}$$

Finally, a simplification removes obviously *true* or *false* parts of the formula.

#### 5. Rejecting Non-Matches

Detecting and rejecting non-matching components as fast and early as possible is probably the single most important step in making deduction-based SCR practical—there

<sup>5</sup>Here,  $\wedge$  means concatenation of lists,  $[]$  the empty list,  $[i]$  a singleton list with item  $i$ , and  $hd$  and  $tl$  the functions head and tail, respectively.

<sup>6</sup>This identification is, however, not always a simple renaming substitution as VDM-SL allows pattern matching and complex data types.

are simply many more non-matching than matching components. Unfortunately, most ATPs are not suited for this task. They exhaustively search for a proof (or refutation) of a conjecture but are practically unable to conclude that it is *not* valid (or contradictory). Therefore, other techniques have to be used to implement rejection filters.

Generally, we may reject a component  $c$  if we find a “counterexample” for its associated proof task  $\mathcal{T}_c$  because it then cannot be valid. Model generators for FOL like Finder [31] or MACE [21] try to find such counterexamples (which are simply interpretations under which  $\mathcal{T}_c$  evaluates to *false*) by systematically checking all possible interpretations. This obviously terminates only if all involved domains are finite, as for example in finite group theory or hardware verification problems. On the other hand, the highly efficient implementation of most model generators (usually using BDD-based Davis-Putnam decision procedures) would make them ideal candidates for fast rejection filters.

However, most domains in our application are not finite but unbounded, e.g., numbers or lists. If we want to use model generation techniques for our purpose, we must map these infinite domains onto finite representations, either by *abstraction* or by *approximation*.

### 5.1. Mapping by Abstraction

One approach to establish this mapping uses techniques from abstract interpretation [6] where the infinite domain is partitioned into a small finite number of sets which are called abstract domains. For each function  $f$  an abstract counterpart  $\bar{f}$  is constructed such that  $f$  and  $\bar{f}$  commute with the abstraction function  $\alpha$  between original and abstract domains, i.e.,  $\alpha(f(x)) = \bar{f}(\alpha(x))$ . E.g., we may partition the domain of integers into three abstract domains  $\{0\}$ ,  $\{x \mid x > 0\}$  and  $\{x \mid x < 0\}$ , called *zero*, *pos* and *neg*. Then, all operations for integers must be abstracted accordingly. For example, for the multiplication  $\times$ , we get the abstract multiplication  $\bar{\times}$  which actually mirrors the “sign rule”:  $neg \bar{\times} pos = pos \bar{\times} neg = neg$ .

Abstract model checking [10] then represents the abstract domains by single model elements and tries to find an abstract countermodel, using an axiomatization of the abstract functions and predicates with a standard FOL model generator. There is, however, a problem. While abstract interpretation may escape to a larger “abstract” domain of truth values in order to make the predicates commute with the abstraction function, standard FOL model generators require the exact concrete domain of *true* and *false* and thus a consistent abstraction may become impossible. E.g., when we try to abstract the ordering on the numbers, *less(zero, pos)* is valid but we cannot assign a single truth value to *less(pos, pos)* because two arbitrary positive numbers may

be ordered either way.

So, while there are some predicates which allow exact abstractions, we have to approximate others. Since we want to use abstract model checking as a rejection filter, we have to make our choices such that the filter produces as few false counterexamples as possible: spurious matches are handled by the subsequent confirmation filter but improperly rejected components are lost forever.

### 5.2. Mapping by Approximation

The second approach to map an infinite domain onto a finite one is done by approximation. From the infinite domain, we select a number of values which seem to be “crucial” for the module’s behavior. E.g., for lists, one usually picks the empty list  $[]$  and small lists with one or two elements (e.g.,  $[a]$ ,  $[a, b]$ ). Then, we search for a model or counterexample. This approach mimicks the manual checking for matches: if one has to find a matching component, one first make checks with the empty list and one or two small lists. If this does not succeed, the component cannot be selected. Otherwise, additional checks have to be applied. This approach, however, is neither sound nor complete. There exist invalid formulas for which a model can be found in a finitely approximated domain (e.g.,  $\forall X : List \cdot \exists i : Item \cdot X = [i]^{\wedge} X$ ), and vice versa (e.g.,  $\exists X, Y, Z : List \cdot X \neq Y \wedge Y \neq Z \wedge Z \neq X$  which has a model only in domains with at least three distinct elements).

While the second case is not too harmful for our application—the performance of the filter just decreases (i.e., more proof tasks can pass), the first one is dangerous: proof tasks describing valid matches might be lost. The experiments which we describe in Section 7.3 are based on this approach. For our prototype implementation we use the model generator MACE [21].

## 6. SETHEO as Confirmation Filter

For the final stage of our filter chain the high-performance theorem prover SETHEO is used. SETHEO is a complete and sound prover for unsorted first-order logic based on the Model Elimination calculus. It accepts formulas in clausal normal form and tries to refute the formula by constructing a closed tableau (a tree of clauses). Completeness is accomplished by limiting the depth of the search space (e.g., with a bound on the size or depth of the tableau) and performing iterative deepening over this bound. In the context of this paper, SETHEO can be seen as a black box which returns “proof found” or “failed to find proof” after the given time-limit. Hence, no further details about SETHEO are given in this paper. For a description of the system and its features see e.g. [16, 24].

With SETHEO’s soundness, we obtain a confirmation filter which guarantees that proof tasks which pass it successfully actually select matching components. Due to our hard time constraints, however, means must be taken not to decrease the recall in an unacceptable way. In the following, we describe how SETHEO has to be adapted in order to be integrated into NORA/HAMMR. We discuss important issues like handling of inductive problems, sorts and equality, and the selection of axioms and parameter settings.

## 6.1. Inductive Problems

Whenever recursive specifications are given or recursively defined data structures are used (e.g., lists) many of the proof tasks can be solved by *induction* only. SETHEO itself cannot handle induction and our severe time-constraints don’t allow us to use an inductive theorem prover. Therefore, we approximate induction by splitting up the problem into several cases. For example, for a query and candidate with the signature  $l : List$ , and the corresponding proof task of the form  $\forall l : List \cdot \mathcal{F}(l)$  we obtain the following cases:  $l = [] \Rightarrow \mathcal{F}(l)$ ,  $\forall i : Item \cdot l = [i] \Rightarrow \mathcal{F}(l)$ , and  $\forall i : Item, l_0 : List \cdot \mathcal{F}(l_0) \wedge l = [i]^{\wedge} l_0 \Rightarrow \mathcal{F}(l)$ .<sup>7</sup> After rewriting the formula accordingly, we get three independent first order proof tasks which then must be processed by SETHEO. This approach can be implemented efficiently. However, we cannot solve every inductive problem.

## 6.2. Equality

All proof tasks heavily rely upon equations. This is due to the VDM-SL specification style and the construction of the proof tasks. While some equations just equate the formal parameters of the query and the library module, others carry information about the modules’ behavior. Therefore, efficient means for handling equalities must be provided. We currently provide two variants: the naïve approach by adding the corresponding axioms of equality (reflexivity, symmetry, transitivity, and substitution axioms), and the compilation approach used within E-SETHEO [24]. Here, symmetry, transitivity and substitution rules are compiled into the terms of the formula such that these axioms need not be added. This transformation, an optimised variant of Brand’s STE modification [3], usually increases the size of the formula, but in many cases the length of the proof and the size of the search space becomes substantially smaller.

<sup>7</sup>Although it would be sufficient to have cases 1 and 3 only, we also generate case 2, since many specifications are valid for non-empty lists only. For those specifications, case 1 would be a trivial proof task which does not contribute to filtering.

## 6.3. Sorts

All proof tasks are sorted. The sorts are imposed from the VDM-SL specifications of the modules and are structured in a hierarchical way. All sorts are static and there is only limited overloading of function symbols. Therefore, the approach to *compile* the sort information into the terms of the formula can be used. Then, the determination of the sort of a term and checking, if the sorts of two terms are compatible is handled by the usual unification. Thus there is no need to modify SETHEO and the overall loss of efficiency is minimal. Our current prototype uses the tool ProSpec (developed within Protein [1]).

## 6.4. Selection of Axioms

Each proof task has to contain—besides the theorem and the hypotheses—the features of each data type (e.g., *List*, *Nat*) as a set of *axioms*. Automated theorem provers, however, are extremely sensitive w.r.t. the number and structure of the axioms added to the formula. Adding a single (unnecessary) axiom can increase the run-time of the prover by magnitudes, thus decreasing recall in an unacceptable way. In general, selecting the optimal subset of axioms is a very hard problem and has not been solved in a satisfactory way yet. Our strong time-constraints furthermore won’t allow us to use time-consuming selection techniques. In our prototype, we therefore use a simple strategy:

1. select only those theories for those data types (e.g., *List*, *Nat*, *Boolean*) occurring in the proof task,
2. within such theories, only select clauses which have function symbols in common with the proof task, and
3. leave out particular clauses and axioms which are known to increase the search space substantially (e.g., long clauses, Non-Horn clauses).

Although this approach is not complete, we use it, since our aim is to solve as many obvious and simple proof tasks (i.e., those which don’t use many axioms or have a complex proof) within short limits of run-time.

## 6.5. Control

Once started, the theorem prover has only a few seconds of run time to search for a proof. This requires that the parameters which control and influence the search (e.g., way of iterative deepening, subgoal reordering) are set in an optimal way for the given proof task. However, such a global setting does not exist for our application domain. In order to obtain optimal efficiency combined with short answer times, *parallel competition* over parameters is used.

The basic ideas has been developed for SiCoTHEO [30] and could be adapted easily: on all available processors (e.g., a network of workstations), a copy of SETHEO is started to process the entire given proof task. On each processor, a different setting of parameters is used. The process which first finds a proof “wins” and aborts the other processes.

## 7. Experimental Results

### 7.1. The Experimental Data Base

All experiments were carried out over a database of 55 list specifications which were modified to have the type *list*  $\rightarrow$  *list* in order to please our still very simple signature matching filter. Approximately 40 of these specifications describe actual list processing functions (e.g., *tail* or *rotate*) while the rest simulates queries. We thus included underdetermined specifications (e.g., the result is an arbitrary front segment of the argument list) as well as specifications which don’t refer to the arguments (e.g., the result is not empty). For simplicity, we formulated the specifications such that the postconditions only used VDM-SL’s built-in sequences.

In order to simulate a realistic number of queries we then cross-matched each specification against the entire library, using partial compatibility as match relation. This yielded a total of 3025 proof tasks where 375 or 12.4% were valid.

### 7.2. Evaluation of Filters

Information retrieval methods [29] are evaluated by the two criteria precision and recall. Both are calculated from the set *REL* of *relevant* components which satisfy the given match relation wrt. to the query and *RET*, the set of *retrieved* components which actually pass the filter. The *precision*  $p$  is defined as the relative number of hits in the response while the *recall*  $r$  measures the system’s relative ability to retrieve relevant components:

$$p = \frac{|REL \cap RET|}{|RET|} \quad r = \frac{|REL \cap RET|}{|REL|}$$

Ideally, both numbers would be 1 (i.e. the system retrieves all and only matching components) but in practice they are antagonistic: a higher precision is usually paid for with a lower recall. We also need some metrics to evaluate the filtering effect. To this end we define the *fallout*

$$f = \frac{|RET \setminus REL|}{|\mathcal{L} \setminus REL|}$$

(where  $\mathcal{L}$  is the entire library) as the fraction of non-matching components which pass the filter as well as the *reduction* which is just the relative number of refuted components. Finally, we define the *relative defect ratio* by

$$dr = \frac{|\mathcal{L} \setminus RET|}{|\mathcal{L}|} \cdot \frac{|\mathcal{L}|}{|REL|}$$

Model	A	B	C
$ List  +  Item $	2+1	3+1	3+2
recall $r$	74.7%	76.5%	81.3%
$\sigma_r$	0.25	0.26	0.25
precision $p$	18.5%	19.6%	16.5%
$\sigma_p$	0.21	0.19	0.16
precision increase	1.5	1.6	1.3
fallout	42.8%	41.0%	55.5%
reduction	50.1%	51.7%	39.0%
defect ratio $dr$	0.51	0.45	0.48

Table 1. Results of model checking

as the relative number of rejected matching components in relation to the precision of the filter’s input. Thus, a relative defect ratio greater than 1 indicates that the filter’s ability to reject only irrelevant components is even worse than a purely random choice.

### 7.3. Rejecting Tasks with Model Generation

For the rejection filter with MACE, we currently use three different models with at most three elements for each data type (in our case *List*, *Item*). Due to the large number of variables in the proof tasks we are generally confined to such small models.

Our experiments with MACE, however, revealed that the restrictions are not too serious. As shown in Table 1, the model checking filter (with a run time limit of 20 seconds) is able to recover at least 75% of the relevant components, regardless of the particular model. The large standard deviation, however, indicates that the filter’s behavior is far from uniform and that it may perform poor for some queries.

Unfortunately, the filter is still too coarse. While each model increases the precision of its answer (compared to the the original 12.4% “precision” of the library) significantly, it still lets too many non-matches pass. The values for fallout indicate that the results in average contain up to 55% of the original non-matching components. Similarly, the overall reduction of approx. 40-50% is at the lower end of our expectations. However, the relative defect ratios show that model checking with any model is at least twice as good as blind guessing.

### 7.4. SETHEO as the Confirmation Filter

For all experiments with SETHEO we used parallel competition with 4 processes exploring different ways of handling equality. Due to technical reasons, we had to restrict the number of modules from our library to 49. This resulted in a total of 2401 proof tasks with 204 or 8.5% matches.

In our first set of experiments we tried to retrieve identical modules from the library (i.e.,  $c \equiv q$ ). The resulting 49 proof tasks are relatively simple and no induction or axioms are needed to prove them. As expected, SETHEO could show all of them within a time-limit of 20 seconds CPU-time (on a sun Ultra-SPARC). Whereas the mean run-time was less than 1s, several proof tasks needed up to 13 seconds.

Then we tried to retrieve matching, but non-identical components. Our experimental basis contains 155 such cases. First, these proof tasks were tried without induction. Here, SETHEO was able to solve 46 proof tasks with a standard set of axioms. The rate of recall could be increased drastically, when our approximation of induction was used. With the same set of axioms, a total of 70 proof tasks could be solved. Due to the increased size of the formulas (esp. in the step case), more overflow errors occurred. Nevertheless, with case splitting we have been able to retrieve 18 matches more than without case splitting. Due to the different structure of the search space, 6 tasks could be shown only without case splitting, making the simple mode interesting for parallel competition. In order to obtain the overall recall of the SETHEO confirmation filter, we have to combine the data of both sets of experiments. From a total of  $204 = 49 + 155$  possible matches, SETHEO could retrieve 125 (49 identical modules, 70 non-identical with case splitting and 6 without case splitting) modules. This yields an overall recall of 61.2%. However, the standard deviation is relatively high as in the model checking experiments, revealing quite different retrieval results for the various queries. Since SETHEO's proof procedure is sound, all solved proof tasks correspond to matches, hence the precision is 100%.

## 8. Related Work

Most early publications on deduction-based SCR (e.g., [20, 28, 14]) were mainly concerned with general conceptual issues and ignored the usability and scaling problems. We will thus discuss only more recent related work.

Zaremski and Wing [34] have investigated specification matching in a slightly more general framework but their main application area is also software reuse. They use the Larch/ML specification language for component description and the associated interactive Larch prover for retrieval. But this promises some severe scaling problems as in our experience only a small fraction of the tasks is provable without interaction. Unfortunately, the paper does not contain any larger experimental evaluation.

Mili et. al. [23] describe a system in which specifications are given as binary relations of legal (input, output)-pairs. They then define a subsumption relation on such pairs and use this for retrieval, relying on Otter to calculate the sub-

sumption. However, their system is still in a prototypical stage, so no relevant statistical evaluation is presented. The examples heavily use auxiliary predicates which are not axiomatized further and thus rely on the arbitrary choice of predicate names to represent domain knowledge. The work of Jeng and Cheng [11] also uses a subsumption test and unfortunately also shares the same problematic confidence in the choice of predicate names. Again, no statistical evaluation is presented.

Scaling problems have been addressed differently. The Inscape/Inquire-system [27] limits the specification language to make retrieval more efficient. Similarly, AMPHION [17] uses a GUI to foster a more uniform specification style which in turn allows an appropriate fine-tuning of the prover. Additional speed-up is achieved by automatically "compiling" axioms into decision theories [18]. These techniques have successfully been applied to assemble over 100 FORTRAN programs from a scientific component library for solar system kinematics. Penix, Baraona and Alexander [26] use "semantic features" (i.e., user-defined abstract predicates which follow from the components' specifications) to classify the components and perform case-based reasoning along this classification to identify the most promising candidates. This classification process uses forward reasoning with an ATP. However, the authors give no evidence of how successful their approach is.

Related work on the use of model checking techniques for infinite domains is much rarer. Jackson [10] is based on [5] and also investigates abstract model checking of software specifications. His goal, however, is to prove the conjectures and not to disprove them. This requires sound approximations which forced him to restrict his logic severely—no negations and exact abstractions only. As soon as approximate abstractions are allowed, this approach also becomes unsound. Wing and Vaziri-Farahani [32] also use abstractions but don't discuss any correctness aspects which are related with them.

## 9. Conclusions and Further Work

In this paper, we have presented NORA/HAMMR, a deduction-based software component retrieval tool. Our goal was to show that such a tool is not only theoretically possible but practical with state-of-the-art theorem provers. We thus designed it as a user-configurable pipeline of different filters. Rejection filters are in charge of reducing the number of non-matching query-component pairs as soon and good as possible. In this paper, we have studied an approach which uses model generation techniques for this purpose. Our experiments with MACE showed that this approach, although neither sound nor complete, returns reasonable results. The final stage of the filter pipeline is al-

ways a confirmation filter which ensures that the selected components really match. Here, we have used the automated theorem prover SETHEO. Even with a short time-limit of 20 seconds, an overall recall of more than 60% was obtained.

We have evaluated our approach with a reasonable large number of experiments. The results obtained so far are very encouraging. We are currently preparing experiments with a library of commercial date and time handling functions as used for example in stock trading software. This work is done in cooperation with the German DG Bank.

Nevertheless, many improvements still have to be made before NORA/HAMMR can really be used in industry. Due to the hard time-constraints (“results while-u-wait”), the reduction of proof-tasks, both in complexity and number is of central importance. Here, powerful rejection filters must ensure that only a few proof tasks remain to be processed by the automated theorem prover. However, our current model-checking filter rejects too much valid matches due to the necessary approximate abstractions. We are thus trying to model exact predicate abstractions with Belnap’s four-valued logic [2] which extends the three-valued LPF consistently. A translation into FOL which reflects the explicit falsehood conditions of Belnap’s logic then yields a sound rejection filter.

Future work will include experiments with specialized decision procedures for the different theories and disproving techniques. Additionally, knowledge-based filters similar to [25] and heuristics will help to reduce the number of tasks to be handled by the confirmation filter. All these filters will be configurable and allow inspection of the behavior of the filter pipeline during each stage of the retrieval.

Current high-performance automated theorem provers are certainly usable as confirmation filters. Much work, however, is still necessary to adapt the ATPs for such kinds of proof tasks. In particular, the requirement of full automatization and the strong time-limits must be obeyed carefully. The experiments showed that parallel competition with several variants (case splitting, set of axioms, handling of equality) is essential to obtain short answer times. Further important issues are the handling of inductive proofs, and the selection of appropriate axioms. Here, powerful heuristics as well as additional information, placed in the data base together with the components (e.g., tactics, lemmas, induction schemes) will be helpful. A further reduction of the search space could be achieved by axiom compilation techniques similar to those of Meta-AMPHION [18]. However, the integration of decision procedures in SETHEO is still an open research topic.

This application of automated theorem proving carries the unique feature that soundness and completeness are not absolutely vital—unsound and incomplete methods only re-

duce the precision and recall of the retrieval tool. This allows interesting and promising deduction techniques (e.g., approximating proofs by filter chains or iteration) to be explored and will help to automate software engineering a little further.

## Acknowledgements

Gregor Snelting originally proposed the filter chain. Christian Lindig helped with the GUI. He and the anonymous referees provided valuable comments. Ralf Haselhorst of the DG Bank, Frankfurt, contributed to the date experiments.

## References

- [1] P. Baumgartner and U. Furbach. PROTEIN: A PROver with a Theory Extension Interface. In [4], pp. 769–773.
- [2] N. D. Belnap. A useful four-valued logic. In *Modern Uses of Multiple-Valued Logic*, pp. 8–37. D. Reidel, Netherlands, 1977.
- [3] D. Brand. Proving theorems with the modification method. *SIAM Journal on Computing*, 4:412–430, 1975.
- [4] A. Bundy, editor. *Proc. 12th CADE*, Vol. 814 of *LNAI*, 1994. Springer.
- [5] E. M. Clarke, O. Grumberg, and D. E. Long. Model checking and abstraction. *ACM TOPLAS*, 16(5):1512–1542, 1994.
- [6] P. M. Cousot and R. Cousot. Abstract Interpretation: A Unified Lattice Model for Static Analysis of Programs by Construction or Approximation of Fixpoints. In *Proc. 4th POPL*, pp. 238–252, 1977. ACM Press.
- [7] J. Dawes. *The VDM-SL Reference Guide*. Pitman, London, 1991.
- [8] B. Fischer. *A systematic approach to type-based software component retrieval*. PhD thesis, TU Braunschweig, 1997. In preparation.
- [9] B. Fischer and G. Snelting. Reuse by contract. In *Proc. ESEC/FSE’97 Workshop on Foundations of Component-Based Software*, 1997.
- [10] D. Jackson. Abstract Model Checking of Infinite Specifications. In *Proc. 2nd FME*, Vol. 873 of *LNCS*, pp. 519–531, 1994. Springer.
- [11] J.-J. Jeng and B. H. C. Cheng. A formal approach to using more general components. In *Proc. 9th KBSE*, pp. 90–97, 1994. IEEE Computer Society Press.
- [12] C. B. Jones and K. Middelburg. A typed logic of partial functions reconstructed classically. *Acta Informatica*, 31(5):399–430, 1994.
- [13] G. E. Kaiser, editor. *Proc. 3rd FSE*, Washington, DC, 1995. ACM Press.
- [14] S. Katz, C. A. Richter, and K. S. The. PARIS: A system for reusing partially interpreted schemas. In *Proc. 9th ICSE*, pp. 377–385, 1987. IEEE Computer Society Press.
- [15] C. W. Krueger. Software reuse. *ACM Computing Surveys*, 24(2):131–183, 1992.
- [16] R. Letz, J. Schumann, S. Bayerl, and W. Bibel. SETHEO: A High-Performance Theorem Prover. *Journal of Automated Reasoning*, 8(2):183–212, 1992.



- [17] M. Lowry, A. Philpot, T. Pressburger, and I. Underwood. AMPHION: automatic programming for scientific subroutine libraries. In *Proc. 8th Intl. Symp. on Methodologies for Intelligent Systems*, Vol. 869 of *LNAI*, pp. 326–335. Springer, 1994.
- [18] M. Lowry and J. Van Baalen. Meta-amphion: Synthesis of efficient domain-specific program synthesis systems. In *Proc. 10th KBSE*, pp. 2–10, 1995. IEEE Computer Society Press.
- [19] Y. S. Maarek and F. A. Smadja. Full text indexing based on lexical relations an application: Software libraries. In *Proc. 12th Annual Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval*, pp. 198–206, 1989.
- [20] P. Manhart and S. Meggendorfer. A knowledge and deduction based software retrieval tool. In *Proc. 4th Intl. Symp. on Artificial Intelligence*, pp. 29–36, 1991.
- [21] W. W. McCune. A Davis-Putnam program and its application to finite first-order model search: Quasigroup existence problems. Technical Report, Argonne National Laboratory, Argonne, IL, USA, 1994.
- [22] K. Middelburg. *Logic and Specification — Extending VDM-SL for advanced formal specification*. Computer Science: Research and Practice. Chapman & Hall, 1993.
- [23] A. Mili, R. Mili, and R. Mittermeir. Storing and retrieving software components: A refinement-based system. In B. Fadini, editor, *Proc. 16th ICSE*, pp. 91–102, 1994. IEEE Computer Society Press.
- [24] M. Moser, O. Ibens, R. Letz, J. Steinbach, C. Goller, J. Schumann, and K. Mayr. The Model Elimination Provers SETHEO and E-SETHEO. *Journal of Automated Reasoning*, 18:237–246, 1997.
- [25] J. Penix and P. Alexander. Toward automated component adaptation. In *Proc. 9th Intl. Conf. on Software Engineering and Knowledge Engineering*, 1997.
- [26] J. Penix, P. Baraona, and P. Alexander. Classification and retrieval of reusable components using semantic features. In *Proc. 10th KBSE*, pp. 131–138, 1995. IEEE Computer Society Press.
- [27] D. E. Perry and S. S. Popovitch. Inquire: Predicate-based use and reuse. In *Proc. 8th KBSE*, pp. 144–151, 1993. IEEE Computer Society Press.
- [28] E. J. Rollins and J. M. Wing. Specifications as search keys for software libraries. In *Proc. 8th Intl. Conf. Symp. Logic Programming*, pp. 173–187, 1991. MIT Press.
- [29] G. Salton and M. J. McGill. *Introduction to Modern Information Retrieval*. McGraw-Hill, New York, 1983.
- [30] J. Schumann. SiCoTHEO — Simple Competitive parallel Theorem Provers based on SETHEO. In *Parallel Processing for Artificial Intelligence 3*, Machine Intelligence and Pattern Recognition 20. Elsevier, 1997.
- [31] J. Slaney. FINDER: Finite domain enumerator. In [4], pp. 798–801.
- [32] J. M. Wing and M. Vaziri-Farahani. Model Checking Software Systems: A Case Study. In [13], pp. 128–139.
- [33] M. Zand and M. Samadzadeh. Software reuse: Current status and trends. *JSS*, 30(3):167–170, Sept. 1995.
- [34] A. M. Zaremski and J. M. Wing. Specification matching of software components. In [13], pp. 6–17.