

Abstract Predicate Entailment over Points-To Heaplets is Syntax Recognition

René Haberland, Kirill Krinkin, Sergey Ivanovskiy[†]
 Saint Petersburg Electrotechnical University "LETI"
 Saint Petersburg, Russia
 haberland1@mail.ru, kirill.krinkin@fruct.org

Abstract—Abstract predicates are considered in this paper as abstraction technique for heap-separated configurations, and as genuine Prolog predicates which are translated straight into a corresponding formal language grammar used as validation scheme for intermediate heap states. The approach presented is rule-based because the abstract predicates are rule-based, the parsing technique can be interpreted as an automated fold/unfold of the corresponding heap graph.

I. INTRODUCTION

In order to understand the term of an *abstract predicate* better let us consider a simple example from *Computational Geometry* [1], the doubly-connected edge list. As the name tells us it is a list of edges, where each edge connects two or three dimensional vertices, usually. If a mesh is triangulated each face consists primarily of three edges, however each vertex coincides twice with a neighbouring vertex, because one edge starts and one edge ends at a vertex. Each vertex may coincide with at least two other vertices. The face structure chosen might improve certain calculations, for instance, the normal vector or seeking for neighbouring faces, etc.

In a points-to heap a location points to some value, in case of the previous example that might be vertices or edges which are pointed to by some variable location. According to the mentioned doubly-linked edge list each list entry may point forward to the next edge or backward to the previous edge. In terms of a points-to heap it is only necessary to specify linked heap entries, non-linked/interleaved memory does not coincide by definition. Imagine, faces, which consist of at least three edges, would need to be specified every time for a heap in a points-to model. That is why a heap predicate may express a complex but intuitive situation very easily, for instance $\text{face}(p1, p2, p3)$ may denote that three vertices $p1, p2, p3$ are connected along a closed circuit building up a face, rather than specifying every time $\exists v1.v2.v3$, s.t. $p1.data \mapsto v1 \star p2.data \mapsto v2 \star p3.data \mapsto v3 \star p1.next \mapsto p2 \star p2.next \mapsto p3 \star p3.next \mapsto p1 \star p1.prev \mapsto p3 \star p3.prev \mapsto p2 \star p2.prev \mapsto p1$.

First of all abstraction means a generalisation, often by introduction of an additional parameters. In this paper the parameter will be primarily Prolog terms, but the underlying assertion language could be abstracted too. By an *abstract predicate* we would like to refer to an arbitrary predicate of the heap (here) in Horn clause form with any number of parameters, which may refer potentially to any number of further abstract predicates. Although introduced by some authors in the

past with capital letters, we tend to characterise those special predicates only with an adjective “*abstract*”, we believe that should be fine. For the example above, $\text{face}(p1, p2, p3)$ stands for the longer \star -separated heaplets. Depending on the actual proof entailment it might be more or less appropriate to perform a fold or an unfold of face causing a proof step to succeed, proceed, fail or block. However, in a fully automated proof, the decision may be hard to find. This paper formalises this problem and chooses a non-traditional approach at a first glance in order to resolve this issue for heaps.

Warren [2] chooses the term “*Programming by Proving*” to stress Prolog can be used as programming language, and what happens, in contrast to other programming languages, is Prolog is searching for a solution – this is what Warren means by proving. Apart from that the *Howard-Curry* isomorphism states there is an interconnection between proving on the one side and programming on the other side. W.r.t. heaplets, the main thesis behind this paper can be summarised in a simplistic way as “*Proving is Parsing*”, meaning by parsing one can prove correctness of specified heap behaviour and that heap behaviour is in a representation close to programming, namely to Prolog rules. A main observation is abstract predicates describe implicitly a language, and so a program’s semantic problem can be resolved by pure syntax recognition.

The structure of this paper is as following: First, current approaches are presented, current restrictions and open questions are provided. Then, heaps and assertions are introduced formally. Abstract predicates and related definitions follow. Afterward the translation problem is formulated. Then the major contribution of this paper is presented – the observation that abstract predicates entailment for points-to heaps can be reduced to syntax parsing, properties of the translation are considered. Finally, implementation details are discussed and conclusions are wrapped up.

II. CURRENT APPROACHES

Since the approach presented in this paper does seem to not follow traditional approaches, only a short review on the closest topics is thought to be meaningful.

First, abstract predicates as presented in this paper are supposed to be understood intentionally close to [3], [4], the most important claim is there is a spatial operator \star which conjuncts two separated heaps and furthermore, for which the rules and axioms of the *Separation Logic* hold which is basically a substructural logic where the contraction rule does not necessarily hold. In this paper abstract predicates just

serve as placeholders for \star -conjoined heap expressions, so no fundamental changes to the core logic are made.

Abstract predicates may be controlled by the user, as it was implemented by `verifast` [5], or semi-automated as done with some tactics applicable to inductively defined predicates as part of `Coq` [6], or fully-automated but with external *fold/unfold* [7] hints for the automated proof entailment. Although the last approach is most convenient from a user's perspective it is challenging from an algorithmic point of view.

Prolog is used in this paper as assertion language. [8] demonstrates how Prolog can be used in order to express in general first-order predicate logic by Horn-clauses. [2] presents what a Prolog semantics in terms of an abstract machine denotes to. Kallmeyer [9] provides a good introduction into Prolog used for parsing. Particularly *adjoint-trees* are proposed as recognition technique of natural language mutations which do not occur in formal or programming languages. However, Prolog is used intensively for implementation. [10] demonstrates by example how Prolog can be used to address ambiguous parsing issues coming from natural language. Matthews discusses recursive tree recognisers in Prolog which seem to match LL(k)-parsers by obeying several modifications to a grammar which is provided as Prolog program: an accepting state needs to be made explicit, rule recursion is simulated by stack which is put onto parameter lists. Matthews uses *difference lists* in order to implement in fact a *partial derivation automaton* for regular languages (see [11]). Finally he proposes Definite Clause Grammars and built-in commands in order to alter Prolog's knowledge base dynamically.

[12] can be considered as a de-facto standard reference on Prolog for natural language processing. From [12] one can find Prolog is incomplete because of possible left-recursive clauses, even when a solution to a given Prolog program exists in general. A model used in order to understand the structure of a natural sentence is a λ -annotated parse tree which – depending on its context – may be interpreted arbitrarily.

III. POINTS-TO HEAPLETS

Definition 3.1. A heap assertion H is inductively defined:

$$H ::= \mathbf{emp} \mid \mathbf{true} \mid \mathbf{false} \mid x \mapsto E \mid H \star H \\ \mid H \wedge H \mid H \vee H \mid \neg H \mid \exists x.H \mid a(\vec{\alpha})$$

The assertion **emp** denotes an empty heap which by default is always true, it is the neutral element for “ \star ” which separates two non-intersecting heaps. The assertion **true** denotes any heap (which always is satisfied), where **false** denotes an arbitrary heap which always interprets as false. This definition is similar to Reynolds' definition [3]. The core component of this definition is $x \mapsto E$, where x is some location identifier (might be an object field, like `o1.field1`), and where E is some valid assignable expression. This definition does no checks on types, this is what is supposed to be processed at an earlier stage [13]. In this paper it is also less of importance if it shall be the immediate meaning or just an address in heap space that stands on the right-hand side (compare with [14]).

Let us consider now, for example, two arbitrary heap predicate definitions formed into Prolog:

$$p1(X, Y) :- \text{pointsto}(\text{loc2}, X), \text{pointsto}(\text{loc3}, Y).$$

$$p2(X, Y) :- \text{pointsto}(\text{loc1}, \text{val1}), p2(X, Y).$$

Here, $p2$ denotes some predicate with two symbols X and Y , which are values pointed by fixed locations `loc2` and `loc3`. In contrast $p1$ is different, hence it refers to predicate $p2$. Whenever we call $p2$ with two syntactically valid term arguments we would have a the form $a(\vec{\alpha})$. Let us remind Prolog just does not find a solution for syntactically valid but semantically invalid terms (meaning the predicate's domain does not include such term).

Interpreting some heap formula H for a given heap maps from a heap domain and a comparison heap into the boolean co-domain, meaning if the provided heap matches the existing heap the interpretation succeeds, and fails otherwise (applying only existing facts and rules which is equivalent to *modus ponens*). The proof in Prolog succeeds iff a true goal is found – a possibility to succeed or a refutation rejects the proof, which with no doubt is exactly the desired behaviour.

For sake of simplicity we agree to canonise heaplets, s.t. we conjoin them all along “ \star ”, having the normal form $a_0 \star a_1 \star \dots \star a_n$ or shorter $\prod_{j=1}^n a_j$ for some $n \geq 0$. Further, \wedge and \vee -connected heaps are turned into Prolog subgoal enumerations of kind $s_j, s_{j+1}, \dots, s_{j+k}$, alternatively split up into separate rule alternatives (or join using the “;”-operator). The negation of an assertion is treated as an functor-guarded negate of a predicate, e.g. `not(P)`, where `not` is a reserved keyword and P some predicate. Just to be mentioned prior to the next sections – negating a sequence of (sub-)goals means those may not follow whilst parsing. The quantification of a fresh variable is allowed spontaneously by simply introducing a fresh variable besides the existing predicate variables.

Moreover, we agree to keep in **emp**, **true** and **false**, although they are syntactic sugar. The heap assertion **emp** is sugar because it could be dropped entirely instead, **true** may be substituted by a tautology in its interpreting conjunction, **false** may be interpreted analogously.

Corollary 3.2. Without going too deep into details, the previous definition implies any valid heap graph can be expressed by induction (and vice versa under the assumption indefinite elements are removed prior to transformation, such as **true** for instance). The proof(s) would be straight forward and inductively defined over the mentioned definition(s).

Definition 3.3. A heap graph is a connected graph within the heap memory section and whose vertices may be pointed by at least one local variable from the stack memory. Each vertex is associated with a heap memory address, its length depends on its data structure. If a vertex is pointing to another vertex, both memory vertices coincide with a directed heap graph edge. If a vertex is pointed by two vertices, then one of which becomes an alias of the other.

IV. ABSTRACT PREDICATES

Prolog is in this paper as programming language in which the assertion language for the heap is specified, and abstract predicates are part of it. Abstract predicates are defined as regular Prolog rules prior to using these rules later in an assertion formula. Often the assertion language does not match a formal specification, particularly for the heap, so this approach is an attempt to bridge this gap by using a logical


```

map([], P, []).
map([X|Xs], P, [Y|Ys]) :-
    Goal =.. [P, X, Y],
    call(Goal), map(Xs, P, Ys).
    
```

Fig. 2. map/3 functional

predicate name t available in Γ .

It is agreed that in a sequence $q_{k,0}, q_{k,1}, \dots, q_{k,m}$ in $q_{m \times n}$ every line is canonised, such that for $s \leq m$ non-trivial entries the first s subgoals are placed and that all remaining $m - s$ subgoals are tautologic subgoals with a domain entirely being true (\top). Moreover, it is agreed that $\exists k.a : -q_k \preceq a : -q_{k+1}$ holds, meaning a predicate rule that occurs earlier in Γ_a has a lower precedence than a predicate that is defined later.

Corollary 4.3. *For the predicate environment Γ of a Prolog program $\Gamma = \bigcup_{t \in T} \Gamma_t$ holds. All Γ_t that depend on each other lay inside a predicate partition $\bar{\Gamma}_t$. $\bar{\Gamma}_t \subseteq \Gamma$ holds.*

Proof: (sketch) The idea behind is to show all dependent $\forall t.\Gamma_t$ lay inside a partition, and all independent partitions, obviously, do not coincide with dependent predicate environments. Naturally, all predicate environments regardless if dependent or independent lay in $\bar{\Gamma}$. Predicates Γ_a and Γ_b from non-coinciding partitions in $\bar{\Gamma}$ can never depend on each other. ■

Remark: Obviously, due to the Halting problem the call of a predicate from a predicate partition may not terminate in general. Next, the expressibility of predicates is considered.

Remark: Possible naming clashes in Γ may be resolved by mangling names including the location of the predicate, such as class where a predicate is defined etc. so the predicates become distinguishable. Predicates within the same location are believed to be together and hence do not clash by definition.

Lemma 4.4. *Abstract predicates cover all first-order predicates.*

Proof: Please refer to [8] for first-order predicate logic completeness expressibility of Prolog. ■

Lemma 4.5. *Abstract predicates may express second- (and even higher-) order logic predicates.*

Proof: Up to this point we were only interested to know through Prolog we can express any first-order predicate logic formula. The following explanation shows we are not restricted ourselves to first-order, but we even can express higher-order in Prolog. In Prolog the built-in predicate `call` allows to call a certain predicate with a list of input and output terms to be passed, for example `pred1(X) :- call(pred2, X)`.

For example, let us define the mapping of a predicate P on an input list, we agree the input parameters shall be encoded in $[X|Xs]$ and output to be $[Y|Ys]$. Then the `map` predicate denotes as shown in Fig. 2. Higher-order predicates may be particularly of interest especially when it comes to dealing with class-objects and *inversion of control*, as it is the case in many *behavioural design patterns*, for instance in the Observer-pattern.

The type of `map/3` is $list_a \rightarrow (list_a \rightarrow list_b) \rightarrow list_b$. So, by introducing third and even higher-order predicates, in analogy to functions we may beat recursion in many cases by using an implicit recursion via higher-order predicates, for instance by application of a left fold that consumes some predicate \oplus and applies it when appropriate having the following signature: `foldl($\oplus :: a \rightarrow b \rightarrow a$, $\varepsilon :: a$, $X :: list_b$) :: a (right folding works in analogy to that). Foldl defines an algebra with an initial value ε and carrier set X and one operation \oplus which is defined on the same type as ε and element-wise for each element of X and calculates a result of same type as ε . For example let us assume we have a equal b are integers and let $X = [1, 2, 3]$ be of kind “list of integers”, let us further agree our initial counter ε equals 7, then foldl will calculate $((\varepsilon + 1) + 2) + 3$ which is 13 which is obviously an integer. ■`

For sake of completeness of the syntax definition from Fig. 1 and the translation in the following section we need to think about how to deal with “;” and “!”. Actually, both are syntactic sugar.

If the body of a predicate rule contains “;” then all right of it has to be split up into a fresh rule with the same name as the origin, so $b : -a_0, a_1, \dots, a_m; a_{m+1}, \dots, a_n$ is split up into $b : -a_0, a_1, \dots, a_m$. and $b : -a_{m+1}, \dots, a_n$. If there is a cut inside $b : -a_0, a_1, \dots, a_m, !, a_{m+1}, \dots, a_n$ then a_0, a_1 until a_m may fail in which case other rules b may be considered as alternatives if any existing. “!” makes sure that if only one subgoal only from a_{m+1} until a_n fails b entirely fails without searching for alternatives. This is again, syntactic sugar, because all possible alternatives may be left-factorised so no other alternatives may be allowed – so, this sugar would insist of rewriting existing predicate rule sets with the same name. This is why without any loss of generality “;” and “!” may in Prolog be dropped from further considerations of expressibility.

Lemma 4.4 and 4.5 conclude that we are able to express all we would like in terms of Prolog, and that we could rewrite some predicate classes without explicit recursion. However, we do not intent to restrict ourselves in terms of μ -recursive predicates.

Definition 4.6. *The predicate unfolding/folding $a(\vec{\alpha})$ of/into some predicate a for some rule system Γ_a with actual term values $\vec{\alpha}$ /subgoals q_k is defined as: (because of lemma 4.5) let Γ_a be w.l.o.g. $a(\vec{y}) : -q_k$ with $q_k = q_{k,0}(\vec{x}_{k,0}), q_{k,1}(\vec{x}_{k,1}), \dots, q_{k,m}(\vec{x}_{k,m})$. Now, if $\vec{\alpha} = (\alpha_0, \alpha_1, \dots, \alpha_A)$ and $\vec{y} = (y_0, y_1, \dots, y_A)$, then $a(\vec{\alpha}) \Leftrightarrow q_{k,0}(\vec{x}_{k,0}), q_{k,1}(\vec{x}_{k,1}), \dots, q_{k,m}(\vec{x}_{k,m})$ with $\alpha_0 \approx y_0, \alpha_1 \approx y_1, \dots, \alpha_A \approx y_A$.*

In case of “ \Rightarrow ” of the above equivalence of $a(\vec{\alpha})$ a predicate is unfold. In case of “ \Leftarrow ” the right-hand side is folded into a predicate `call`. “ \approx ” stands for term unification.

V. INTERPRETING ABSTRACT PREDICATES OVER HEAP AS SYNTAX RECOGNITION

The goal of syntax recognition is to automate the check of heap predicates in specifications against an inferred heap state. The technique applied is syntactic for a semantic problem. The problem with predicates is the non-determinism of

when to fold/unfold. Proof tactics have been implemented in theorem provers and checkers, like Coq [6], in very dedicated domains only, but the quality of the fold/unfold is still far from satisfiable. An automated fold/unfold approach would be highly desirable, so additional specifications or even manual interactions can be zeroed. The new approach presented in this paper requires the following steps:

- 1) Translate incoming program and annotated assertions into Prolog terms which are integrated into [13].
- 2) Define abstract predicates within an annotated section in the incoming program. These Prolog rules need to be syntactically correct.
- 3) Generate formal grammar for found abstract predicates. File grammar over to a parser-generator which will finally emit a valid and concrete parser.
- 4) While running the proof, trigger certain parse rules depending on abstract predicate calls found.

It will provide a different non-standard way of dealing with the problem.

Observation 5.1. *When looking at how a heap is generated it reminds a production system for formal languages.*

Terminals become points-to expressions (cmp. with definition 4.2) or relations, and non-terminals become abstract predicate subgoals. Terminals are concatenated, where points-to expressions are loosely coupled with \star . \star is commutative. Points-to expressions may be concatenated too, when the pairs are ordered according to the left-hand side location name. If local names clash, a namespace would resolve this by full qualification, e.g. by name prefixes.

Theorem 5.2. *The predicate partition builds up an production system and manifests in fact a context-free grammar.*

Proof: The left-hand side of a predicate rule may only be no more than one non-terminal. It is more than regular because there is no such claim the right-hand side needs to be right-recursive. If the head of the rule contains arguments this still does not change statically the dependencies in between the predicates. A predicate partition has one starting non-terminal. ■

Observation 5.3. *When looking at how a heap is derived from abstract predicates one may think about reducing it.*

The implication underneath, however, would be both, from satisfiable. An automated fold/unfold approach would be from satisfiable. An automated fold/unfold approach would be unfolded until both sides establish an equivalence. This would be a bi-directional approach. However, that problem could be reduced to *Post's Correspondence Problem* and is unfortunately undecidable in general, hence is not considered here any further.

Observation 5.1 seems promising, so the heap predicate check can be re-formulated as: “Given a \star -connected heap, does it match a given heap specification or not?”. But there might be further questions related, such as: “What would be the closest correct heap, s.t. it satisfies the current heap specification?”, which could deliver us answer to the counter-example problem.

Lemma 5.4. *The word-problem for abstract predicates P can be formulated as: Given $\alpha_1, \alpha_2 \in L(G(P))$ does $\alpha_1 \equiv \alpha_2$ hold? $G(P)$ denotes the formal context-free grammar obtained from the predicate partition of P .*

Proof: Here $\alpha = (a + A)^*$, and $a \in T$, $A \in NT$. T denotes all terminals which are parameterised and encode source and target of “ \mapsto ”, NT denotes non-terminals which contain all predicates and parameterise all valid input terms. The rule set P is the translated set of predicates, $L(G(P))$ is the language generated by the generated grammar by Prolog rules. The starting non-terminal is a predicate call from either α_1 or α_2 . Regardless of the particular kind of parser to be used the follow set $\sigma(\alpha)$ and the first-terminal sets $\pi(\alpha)$ need to be calculated (see later). One important fact is there is not a single start non-terminal, but there might be several depending on number of predicate calls in α_1 and α_2 . Furthermore, there is not only one path searched from $\alpha_1 \vdash^* \alpha_2$ but also $\alpha_2 \vdash^* \alpha_1$. Only if there is no path found in both directions α_1 does not coincide with α_2 , otherwise it does. In order to check if two sentences match, it is not only necessary to construct paths between predicates, it is also necessary to consume initial and intermittent `pointsto`-terminals. Parameters on terminals and non-terminals shall be bound according to the current binding and unified with α_1 and α_2 . ■

VI. TRANSLATION OF HORN-CLAUSES INTO GRAMMAR

This section considers how abstract predicates provided as Prolog rules are translated into a general context-free grammar. Before that, Prolog rules need to be analysed lexically, so all expressions of form $loc \mapsto val$ are considered tokens. Multi-paradigmatic programming [16] allows interpreting Prolog rules during execution of some main Prolog application, which, from satisfiable. An automated fold/unfold approach would be needed to be done once until all abstract predicates have been processed. The translation process from Prolog rules into a formal grammar is astonishing simple. However, Prolog rules may have argument terms on the left and the right side of “ $:-$ ”, this can be modelled by introducing attributes to the generated grammar, we finally obtain an attributed grammar. Hence, the translation $C \llbracket \cdot \rrbracket$ can be defined straight:

Definition 6.1. $C \llbracket \cdot \rrbracket$ is defined as rule transducer for an incoming abstract predicate set and attributed grammar as output:

$$\begin{aligned} C \llbracket \cdot \rrbracket &= \emptyset \\ C \llbracket C_1.C_2 \rrbracket &= C \llbracket C_1 \rrbracket \cup C \llbracket C_2 \rrbracket \\ C \llbracket a(\vec{x}) : -q^0(\vec{x}), \dots, q^n(\vec{x}) \rrbracket &= \{a_{\vec{x}} \rightarrow q_{\vec{x}}^0 \dots q_{\vec{x}}^n\} \end{aligned}$$

In contrast to previous notations subgoals here have upper indices and \vec{x} now accommodates all variable symbols within of a predicate rule for notational comfort, so if a particular subgoal q^j for some j does not require all components of \vec{x} then it does not. Remind \cup is a set union where the element sequence matters. As can be seen from both notations are pretty similar to each other and are interchangeable, the inverse operation $C^{-1} \llbracket \cdot \rrbracket$ translates an attributed grammar back into Prolog and can be defined as:

Definition 6.2. $C^{-1} \llbracket \cdot \rrbracket$ is defined as rule transducer for an incoming attributed grammar and an abstract predicate set

as output:

$$\begin{aligned} C^{-1}[\Box] &= \emptyset \\ C^{-1}[[C_1 \ C_2]] &= C^{-1}[[C_1]] \cdot C^{-1}[[C_2]] \\ C^{-1}[[a_{\vec{x}} \rightarrow q_{\vec{x}}^0 \dots q_{\vec{x}}^n]] &= \{a(\vec{x}) : -q_0(\vec{x}), \dots, q_n(\vec{x})\} \end{aligned}$$

Corollary 6.3. $C[\Box]$ and $C^{-1}[\Box]$ terminate for any well-defined domain input.

Proof: The proof is rather trivial, since there is no infinite cycle possible. Both, $C[\Box]$ and $C^{-1}[\Box]$ linearly scan all incoming rules successively from left to right. Suppose, there was a cycle in between particular rules. Even so, both translations will finally terminate because cycles may bother only while parsing, not while translating. The starting point to a predicate partition corresponds one to one to the starting non-terminal of a subgrammar. There might be several entry points for abstract predicates, and so are the entry points corresponding to non-terminals for a grammar. ■

We still need to investigate $C[\Box]$ and $C^{-1}[\Box]$ soundness and completeness.

Corollary 6.4. $C[\Box]$ and $C^{-1}[\Box]$ are total and both mappings are complete and sound.

Proof: It is not hard to verify that $C \circ C^{-1} \circ C \equiv C$ hold as well as $C^{-1} \circ C \circ C^{-1} \equiv C^{-1}$ by simple substitution of the definitions from above. Because of the discussions in section III, “!” nor “;” do not matter w.r.t. expressibility. If, however, the domain of $C[\Box]$ is supposed to not terminate, then its co-domain will cause exactly the same behaviour, same holds for $C^{-1}[\Box]$. ■

VII. PARSING

For the purpose of a simple and intuitive algorithm the constants from definition 3.1 will not be considered. Because as mentioned earlier they are not intrinsic and can be dropped therefore. Essentially those heap constants provide partial heap expressiveness and may be considered for future research, w.r.t. class objects a **true** could possibly mean, for instance, to consume all `pointsto` until a (rule-dependent) marker pops up indicating a *safe* synchronisation point in terms of error productions [17] for the ongoing parsing as described briefly. The input word is finite, however in general the number of unfolds may be hypothetically infinitely many – but shown later the π -function allows us to pre-calculate the following terminal symbols with polynomial efforts.

This section introduces fundamental conventions necessary to complete some generic LL(k)-parser for demonstration purposes. Fundamentally, this is not only for an unlimited forward-looking LL parser, it is also possible to use some different parser, for instance based on a generalised LR or Earley parser. First, a sentence is defined as some composite of `pointsto` (terminals) and further subgoals with term arguments (non-terminals) – something that the right-hand side of a (Prolog) rule comes up with. Second and third, in analogy to a LL(k)-parser but with functional space rather than single character as for strings the definitions of *first* and *follow* sets are introduced. Forth, both SHIFT and REDUCE are proposed for some general parser implementation.

Definition 7.1. An abstract sentence α is a \star -conjunction of heaps which are denoted by $a \mapsto b$, where a is some unique location identifier and b some value object or **nil**.

For example, $\alpha ::= [\text{pointsto}(x, \text{nil}), \text{pointsto}(y, 1), \text{member}(x, [y])]$ may describe the current state of the heap during the verification of an imperative program, and \star is replaced in the previous list by commas. The specification of a rule may insist on $[\text{pointsto}(Y, 1), \text{member}(X, [Y|_])]$, $\text{pointsto}(X, _)$. So, what is necessary to check the abstract sentence from the program matches the sentence from the specification is primarily to check whether all parts from either of both is derivable from each other.

An abstract sentence may also contain term unification, such as `pointsto(X, 5), X=Y`. Term unification has to be thoroughly analysed and separated from the remaining two cases, namely `pointsto` which denotes terminals, and predicate calls as subgoals which denote non-terminals later on. It has to be taken into consideration that not any terms may be unified, since there is a occurs-check taking place by intention w.r.t. the given implementation and by default in Prolog, so indefinite terms or recursive term definitions are strictly prohibited.

Let us now formulate an initial algorithm in order to check two abstract sentences match or do not match, let us consider algorithms 1. π denotes the function being introduced shortly. The problem is we reduce (factual predicate unfolding) possibly ad absurdum, we do not know determined when and how often to fold and unfold which obviously also depends on the rules themselves. Assume, we had some $\alpha_1 = [\underbrace{a \mapsto b}_{i_0}, \dots, i_{m_1}, \underbrace{q_1(x)}_{p_0}]$ and some

$\alpha_2 = [\dots, \underbrace{a \mapsto b}_{j_3}, \dots, \underbrace{q_1(x)}_{q_7}]$ we would like to match against

with. SHIFT-TERMS would first of all unify i_0 with j_3 and possibly continue with all other matching terms. REDUCE-PREDS will try match all matching predicates which may have to be unfolded first, that is why the first terminal of a predicate may be required first, the expansion $\text{expand}(p_k, \alpha_1)$ expands the predicate head by the corresponding body of p_k (compare with definition 4.6 and Fig. 1) into a new abstract sentence α' which might be described in Prolog as $\text{concat}(\alpha, [i_7, i_8, i_9], \alpha')$, if for instance $q_1(x)$ unfolds into $[i_7, i_8, i_9]$.

Definition 7.2. The first set is defined as co-domain of a total map π with type $(T \cup NT) \rightarrow 2^T$ for $m > 0$, such that:

$$\pi(a) ::= \begin{cases} a & \text{if } a \text{ is } X \mapsto Y \text{ or } \Gamma_a ::= a. \\ \bigcup_{0 \leq j \leq n} \pi(q_{j,0}) & \text{if } \Gamma_a ::= a : -q_{m \times n}, n > 0. \end{cases}$$

Essentially, π determines all terminals that start with `pointsto` or are beginnings (only the first terminal) of predicate subgoals (independent of its arguments).

Definition 7.3. The follow set $\sigma(t) \subseteq T$ for $t \in (T \cup NT)$ is defined as:

Algorithm 1 A naïve algorithm for word equality check for abstract sentences; **Input:** $\alpha_1 = [i_0, \dots, i_{m_1}, p_0, \dots, p_{n_1}]$, $\alpha_2 = [j_0, \dots, j_{m_2}, q_0, \dots, q_{n_2}]$ with i and j as `pointsto`-terminals, and subgoals p and q as non-terminals. Γ contains all predicate definitions. **Result:** *True* in case α_1 equals α_2 , *false* otherwise.

```

1: procedure SHIFT-TERMS( $\Gamma, \alpha_1, \alpha_2$ )
2:   for all  $k \in \{0, \dots, m_1\}$  do
3:     if  $\exists l.l \in \{0, \dots, m_2\} \wedge i_k \approx j_l$  then
4:        $\alpha_1 \leftarrow \alpha_1 \setminus i_k$ 
5:        $\alpha_2 \leftarrow \alpha_2 \setminus j_l$ 
6:     end if
7:   end for
8: end procedure

9: procedure REDUCE-PREDS( $\Gamma, \alpha_1, \alpha_2$ )
10:  for all  $k \in [0..n_1]$  do  $\triangleright$  try match with terminal
11:    if  $\exists l.l \in [0..m_2] \wedge j_l \in \pi(p_k)$  then
12:       $expand(p_k, \alpha_1)$ 
13:    else  $\triangleright$  try to reduce in  $\alpha_2$ 
14:      if  $\exists l.l \in [0..n_2] \wedge (\pi(q_l) \cap \pi(p_k) \neq \emptyset)$  then
15:         $\alpha_1 \leftarrow expand(p_k, \alpha_1)$ 
16:         $\alpha_2 \leftarrow expand(q_l, \alpha_2)$ 
17:      else  $\triangleright$  Match
18:        if  $m_1 = m_2 = n_1 = n_2 = 0$  then
19:           $true \rightarrow Halt!$ 
20:        else  $\triangleright$  Non-Match
21:           $false \rightarrow Halt!$ 
22:        end if
23:      end if
24:    end for
25: end procedure
    
```

$$\sigma(t) ::= \begin{cases} \bigcup_{i,j} \pi(q_{i,j+1}) & \text{if } t \text{ is at } (i, j < n) \text{ in } q_{m \times n} \\ & \wedge 0 \leq i \leq m \\ & \wedge q_{i,j+1} \neq \top \\ & \wedge \exists a.\Gamma_a ::= a : -q_{m \times n} \\ \bigcup_a \sigma(a) & \text{if } t \text{ is at } (i, n) \text{ in } q_{m \times n} \\ & \wedge \Gamma_a ::= a : -q_{m \times n} \\ & \wedge \exists b.\Gamma_b ::= b : -q_{m_b \times n_b} \\ & \wedge a \text{ is at } (i_b, j_b) \text{ in } q_{m_b \times n_b} \\ \emptyset & \text{otherwise} \end{cases}$$

The follow set determines literally all terminals that may follow a `pointsto` or a predicate subgoal from all considered rules. Now we have defined π and σ we are able to synthesise from this a LL(k) recogniser ([17] might be a helpful introduction).

Example 7.4. Given the following production rules $q_1 \rightarrow a$, $q_2 \rightarrow aq_2$ | q_3b , $q_3 \rightarrow \varepsilon$ | q_3a these rules are obviously ambiguous, for instance in $\pi(q_2) = \{a\}$, $\sigma(a) = \{\varepsilon\} \cup \pi(q_2) \cup \pi(q_3) \cup \sigma(q_3)$.

Example 7.5. Given the following finite specification $[(loc1, v1), p1(loc1, loc2), (loc2, v2)]$ and the word $[(loc1, v1), (loc2, v2)]$ which is true, but only if $p1$ does not dump a heap.

In case the input word is not particularly a sequence of terminals, but an abstract sentence, it will be necessary to cut redundant calculations as early as possible. Hence, memoizing calculated subgoals would not only enhance the speed of search (since only the predicate name and its parameters play a role in memoization), it would resolve matching the first non-terminal issue, which by the way, matches neatly in the described definition of π and σ .

Negated predicates are dropped here, refer to section IX for details.

VIII. PROPERTIES

In Fig. 3 a sample heap configuration is shown. This configuration consists of 8 triangles, where each triangle is crossed by either a dotted or a tortuous line. The line from the midpoint of v_0 and v_3 to M_1 denotes the triangle $\Delta(v_0, v_3, M_1)$, where the tortuous line between the midpoint of v_0 and v_1 and M_1 addresses $\Delta(v_0, M_1, v_1)$. And so, Fig. 3 demonstrates there might be more than one way of spanning the heap graph by some provided heap predicate set, namely either by the triangles marked with dotted lines or with the tortuous lines. It is, however, essential that all vertices need to be included in a heap in order to decide heap graph isomorphism.

There is (currently) only one position where a non-deterministic decision has to be taken out, namely the decision where to bound the input stream w.r.t. object boundaries. If introducing a convention that only common fields of an object are allowed and are canonised, the decision becomes determined. Obeying the convention could be performed within polynomial effort, as the rest of the parsing routine – which is well-described and is known to be tractable within polynomial efforts (refer to [17] on parsing foundations). If all predicate partitions are parsable, as described with possible practical restrictions discussed in the previous sections, then entailment over points-to heaplets becomes decidable and finally terminates with an answer or proof refutation. The proof refutation will be in fact be a syntax error with corresponding coordinates in the points-to encoded input word according to the expected predicate partition.

The core memory model has not been modified nor extended, except the introduction of abstract predicate definitions over the existing points-to model. The proposed extension shall therefore be compatible e.g. with Reynolds' [3] or Burstall's model [14], however in this notation in contrast to Burstall cell addresses where not used, so this approach is conventionally

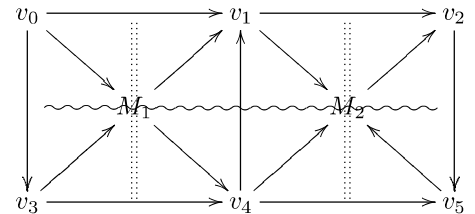


Fig. 3. A sample heap configuration where v_j with $j \neq 2, j \neq 5$ outgoing edges are of class object type with more than one associated pointer attributes

closer to Reynolds. The consequences of Burstall’s notation could be researched, since it is common practice compound objects are referred in practice by reference, not by its content.

IX. IMPLEMENTATION

The implementation is in GNU Prolog and uses ANTLR version 4 [18], and is supposed to incorporate the framework presented in [13], [19]. Initially both tools were chosen for simplicity and extensibility reasons and for lecturing purposes.

In order to gain from flexibility and a huge support of existing software packages, the chosen programming paradigm is *multi-paradigmatic* [16], which allows the developer to write and run programs in different programming languages and profit from both of its advantages. The integration of both works astonishingly simple due to a Proxy design pattern and an interpretation in both directions. There exists also an experimental user interface based on tuProlog for prototyped development. The implementation first translates input language into an intermediate representation which are Prolog terms, afterwards the assertions are copied separately into a Prolog theory, and abstract predicates are transformed into a ANTLR 4 grammar file as explained earlier. Whenever a parsing is requested, an internal syntax recognition process is initialised in the language the ANTLR outputted the recognisers (which is Java here). The output and control is passed back to the invoker. This way abstract predicates can be checked fully automated, and in case of an error the corresponding error will be processed.

ANTLR makes use of so-called “*syntactic and semantic predicates*” in order to fight syntax ambiguity. Since ANTLR does not necessarily cover in practice all LL(k)-grammars strictly, there is of course still room for improvement. Practically this means that occasionally there may appear grammars which shall, but which do not parse due to current limitations of the ANTLR parser generator. There were made experiences with other parser generators, e.g. LR(0)-parsers resolve this issue and even recognise left-recursion by definition, but lack from known shift-reduce restrictions on the other side therefore, like with GNU bison, for instance. A good introduction to parsing techniques can be found in [17].

As an example of required transformation during the analysis of lexemes and tokens some precautions were required. First, $bar \mapsto foo$ is mangled to pt_3bar_3foo where 3 is the length of the name or corresponding accessor. If the accessor is compound, e.g. $b.f.g$, then the accessor length avoids ambiguity. For instance $pointsto(X, 2)$ is mangled to a Prolog atom p_X_2 . If needed, a mangled name can be demangled – the internal parsing may be done by a Java helper which is then visible in Prolog as a helper built-in predicate via [16]. Second, the left-hand side (de-)canonisation (on Prolog level). $p1(X, [X|Y]) :- \dots$ is transformed into $p1(X1, X2) :- X1=X, X2=[X|Y], \dots$. Third, a Prolog rule $p(X, Y) :- \alpha$ may be translated to a ANTLR-grammar snippet: $p[String X, String Y] : \alpha$. This way all synthesised attributes may be passed top-down, inherited attributes may be modified to some predicate p by adding `returns` together with the inherited attribute names just before colon. So, all that is necessary is to decide whether a variable is inherited or synthesised in order to decide its position in the grammar snippet.

When abstract predicates are turned into a concrete grammar, e.g. a ANTLR grammar file, the problem arises that unified terms are together with `pointsto` terminals and non-terminals. Unifying terms need to be separated from terminals and non-terminals, therefore they are moved into translating rules beside the attributed grammar. For instance, ANTLR introduces translating rules using the curly brackets within a sentence. Negated sentences and fragments of it can be introduced by “ \sim ” and brackets, and mean the included sentence may not appear. No further cases are discussed since either by attributes or translating rules additional behaviour may be mimicked, such as a failing predicate as a parser signal.

X. CONCLUSIONS

The approach presented proposes a technique to automatically entail points-to heaplets by syntax analysis rather than manually fold/unfold abstract predicates. If a predicate partition is representable as a valid set of a parser’s rule set then there will be a definite answer whether heap specification and an existing heap configuration match. The model used in between Prolog and a concrete parser rule set is an attributed grammar [17] which translates inheriting and synthesising attributes which correspond to Prolog’s head terms. We believe the implementation of stack frames which is different to those of common imperative programming language, gives Prolog an important advantage in reasoning since it is what we would expect from the attribute content control – without the need of additional implementation nor development costs, because it is part of Prolog’s core [2]. The used points-to heaplets may correspond to a Separation Logic styled model obeying its axioms and rules. It is true, Prolog’s abstract machine is an interpreter and therefore on average slower than any natively compiled code, but the question of performance is minor in this case since we are mostly interested in exploring tractability and expressibility first of all – since verification is done separately from the generated code, we take intentionally the risk of being a bit slower occasionally.

Instead of *simulating* only symbols within an artificially introduced assertion language, the assertions here are all expressed in the same language in which proofs will be taken out, Prolog. There is no overhead of mapping in between assertion and proof language as it is for instance the case in [6]. Symbols may be used very closely to the first-order predicate logic, symbols and terms may be unified, which is a bi-directional reasoning technique. It is believed expressibility in Prolog terms, especially with regards to objects, may improve over non-adequate representations, as it was proven concept in [20] on markup-notations for terms. It remains an open question due to term unification and the rule-based predicate partitions if and how error production rules may in fact advance the completeness of reasoning rules further, for instance w.r.t. abduction and proof explanation. Another practical advantage using Prolog whilst proving is the possibility to load parts or none of Prolog rules and facts in order to try some question without loading all at once, which makes error tracing comfortable.

Open questions and future work includes the possibility of partial heap specifications using constant keywords like **emp**, **true** and **false**, and the question if a proof may get simpler just if the heap graph is required to be connected. The last question

arises for improved error location on proof refutations using error production rules to be introduced and invoked during parsing. It could stop, for instance, on a refutation and behave like in a “panic mode” [17] consuming all input tokens until a synchronised save state or a sequence of consecutive known safe tokens.

ACKNOWLEDGEMENT

†This article is dedicated to Sergey Ivanovskiy in remembrance.

Parts of this paper are prepared as a contribution to the state project of the Board of the Ministry of Education of the Russian Federation (task # 2.136.2014/K).

REFERENCES

- [1] M. de Berg, O. Cheong, M. van Kreveld, and M. Overmars, *Computational Geometry: Algorithms and Applications*. Santa Clara, California, USA: Springer-Verlag Berlin Heidelberg, 3rd ed., 2008.
- [2] D. H. Warren, “Applied logic - its use and implementation as a programming tool,” Tech. Rep. 290, SRI International, 333 Ravenswood Ave, Menlo Park, CA 94025, USA, June 1983.
- [3] J. C. Reynolds, “Separation logic: A logic for shared mutable data structures,” in *Proceedings of the 17th Annual IEEE Symposium on Logic in Computer Science*, (Washington, DC, USA), pp. 55–74, IEEE Computer Society, 2002.
- [4] M. Parkinson, *Local Reasoning for Java*. PhD thesis, Cambridge University, 2005.
- [5] B. Jacobs, J. Smans, P. Philippaerts, F. Vogels, W. Penninckx, and F. Piessens, “Verifast: A powerful, sound, predictable, fast verifier for C and Java,” in *Proceedings of the Third International Conference on NASA Formal Methods, NFM’11*, (Berlin, Heidelberg), pp. 41–55, Springer-Verlag, 2011.
- [6] Y. Bertot, P. Castran, G. Huet, and C. Paulin-Mohring, *Interactive Theorem Proving and Program Development - Coq’Art: The Calculus of Inductive Constructions*. Texts in theoretical computer science, Berlin, New York: Springer, 2004.
- [7] G. Hutton, “Fold and unfold for program semantics,” in *In Proc. 3rd ACM SIGPLAN International Conference on Functional Programming*, (Baltimore, Maryland), pp. 280–288, ACM Press, 1998.
- [8] R. A. Kowalski, “Predicate logic as programming language,” in *Information Processing (IFIP)*, pp. 569–574, North-Holland Publishing, 1974.
- [9] L. Kallmeyer, *Parsing Beyond Context-Free Grammars*. Cognitive Technologies, Springer Heidelberg, 2010.
- [10] C. Matthews, *An Introduction to Natural Language Processing Through Prolog*. White Plains, NY, USA: Longman Publishing, 1st ed., 1998.
- [11] J. A. Brzozowski, “Derivatives of regular expressions,” *Journal of the ACM*, vol. 11, pp. 481–494, Oct. 1964.
- [12] F. Pereira and S. Shieber, *Prolog and Natural-Language Analysis*. Microtome Publishing, 2002.
- [13] R. Haberland and S. Ivanovskiy, “Dynamically allocated memory verification in object-oriented programs using prolog,” in *Proc. of the 8th Spring/Summer Young Researchers’ Colloquium on Software Engineering (SYRCoSE 2014)* (A. Kamkin, A. Petrenko, and A. Terekhov, eds.), pp. 46–50, 2014.
- [14] R. M. Burstall, “Some techniques for proving correctness of programs which alter data structures,” in *Machine Intelligence* (B. Meltzer and D. Michie, eds.), vol. 7, pp. 23–50, Scotland: Edinburgh University Press, 1972.
- [15] J. Berdine, C. Calcagno, and P. W. O’Hearn, “Symbolic execution with separation logic,” in *Programming Languages and Systems, 3rd Asian Symposium*, (Tsukuba, Japan), pp. 52–68, November 2005.
- [16] E. Denti, A. Omicini, and A. Ricci, “Multi-paradigm Java-Prolog integration in tuProlog,” *Science of Computer Programming, Elsevier Science*, vol. 57, pp. 217–250, Aug. 2005.
- [17] D. Grune and C. J. H. Jacobs, *Parsing Techniques: A Practical Guide*. Upper Saddle River, New Jersey, USA: Ellis Horwood publ., 1990.
- [18] T. Parr, *The Definitive ANTLR 4 Reference*. The Pragmatic Programmers, 2012.
- [19] R. Haberland and K. Krinkin, “Verifikaciya ob’ektno-orientirovannyx program s pomo’shu prologa (russian),” in *Izvestiya LETI*, vol. 11, 2015.
- [20] R. Haberland and I. L. Bratchikov, “Transformation of XML documents with Prolog,” in *Advances in Methods of Information and Communication Technology*, vol. 10, pp. 99–111, 2008.