

Deductive Verification of Sparse Sets in Why3

Catherine Dubois

ENSIIE, INRIA, Université Paris-Saclay, LMF, France
catherine.dubois@ensiiie.fr

Abstract. To represent finite sets of integers on an interval $0..n$, Briggs and Torczon studied a very simple data structure in 1993, called *sparse sets*. With this representation, initialization, membership test, insertion and deletion of an element are $O(1)$ operations. This data structure is often used in compilers to allocate registers or to represent the objects in a video game. A variant of this data structure is also used in finite domain constraint solvers to represent the domains of integer variables. This variant makes it a backtrackable data structure. We have formalized and verified the original data structure and its variant in Why3. Set operations such as intersection and union are formally verified, even though they are less commonly used with this representation of sets. To our knowledge this is the first formal verification of the backtrackable variant of sparse sets used as domains.

1 Introduction

Sets are seldom primitive objects in programming languages or specification formalisms. There are such objects in the old programming language Set1 [21] or in the logic programming language $\{\text{log}\}$ [8] and also in the formal languages B [1] or Event-B [2] and TLA+ [14]. More usually, they are available as implementations in libraries, based on underlying data structures such as sorted lists, red-black trees, AVL trees, B trees, skiplists, etc. In this paper, we focus on sparse sets, studied by Briggs and Torczon [4] in 1993, also appearing as an exercise in the famous book "The Design and Analysis of Computer Algorithms" written by Aho and Hopcroft [3]. This data structure dates back to computer folklore, it is used in different applications like register allocation, video game, constraint solving. With this mutable representation based on arrays and simple manipulations, initialization, membership test, insertion and deletion of an element are $O(1)$ operations. Many implementations exist on the web, in several languages, e.g. Java, C++, C, Rust.

Sparse sets appear as a benchmark (Constant-time sparse array) of VACID-0, a suite of benchmark verification problems proposed in 2010 [17]. The sparse sets data structure as it is described by Briggs and Torczon is a particular case of the latter in which there is one less indirection. A solution¹ where 3 operations

¹ available at https://toccata.gitlabpages.inria.fr/toccata/gallery/vacid_0_sparse_array.en.html

(membership, add and remove) are implemented, has been given using Why3 by Filliâtre and Paskevich.

Sparse sets are also used in constraint solvers as an alternative to range sequences or bit vectors for implementing domains of integer variables [15] which are nothing else than mathematical finite sets of integers. Sparse sets as domains are slightly different from sparse sets introduced by Briggs and Torczon making them very easy to store and restore when backtracking for finding solutions.

Our main contribution is a formally verified implementation of sparse sets as domains and its various operations, developed with the deductive verification tool Why3 [12], extracted in OCaml. In addition to classical set operations (test membership, remove, etc.), we specify and verify an operation that allows the user to undo some operations very easily (in one simple assignment). This contribution brings some more confidence in the data structures used in constraint solvers, as it has been done by Ledein and Dubois [16] for the traditional implementation of domains as range sequences.

In [7], Cristiá and the author have formalized this sparse set as domain variant and verified three simple operations (remove, bind and membership) in three formalisms Why3, EventB and $\{\log\}$. However they do not address the verification of the backtrackable dimension and in particular the undo operation.

The article is structured as follows. In Section 2 we give an informal overview of sparse sets. In Section 3, we briefly introduce Why3 and WhyML. In Section 4, we detail our WhyML implementation of sparse sets and discuss its deductive verification. In Section 5, we first introduce the modifications to the data structure when it is used to represent the domain of integer variables in constraint solvers, then we present the WhyML formalization of this backtrakable variant by focusing mainly on the additional artefacts we used to verify the undo operation. Section 6 presents some tests performed on the OCaml code extracted from our models. Finally we conclude and present some future work.

All the code described in this paper is available on https://gitlab.com/cdubois/why3_sparsesets.

2 Sparse sets

Sparse sets are used to represent subsets of natural numbers up to $N - 1$, where N is any non-zero natural number. The range $0..N - 1$ is called the universe of the sparse set in the following. A sparse set D is represented by two arrays of length N called *Dense* and *Sparse*, and a natural number $sizeD$ ². The current elements of the finite set are those in $Dense[0, sizeD - 1]$ — let us call this subarray the effective part —, the rest of the array being *garbage*. The array *Sparse* maps any value $v \in [0, N - 1]$ to an index ind_v in *Dense* or is not initialized. Thus, for the current elements v of the set, $Sparse[v]$ has a value i in the range $[0, sizeD - 1]$ and $Dense[i]$ is equal to v . If D is empty (resp. the full set), $sizeD$ is equal to 0 (resp. N).

² The name of this data structure may be explained by the fact that the *Sparse* array may have holes whereas the *Dense* array is more compact.

The two invariants of the data structure representing the set D are as follows:

$$D \subseteq 0..N-1 \wedge D = \{Dense[i] \mid 0 \leq i < sizeD\} \quad (P_1)$$

$$v \in D \iff 0 \leq Sparse[v] < sizeD \wedge Dense[Sparse[v]] = v \quad (P_2)$$

Fig. 1a illustrates this representation. This state has been reached after inserting the elements 3, 6, 4, 7, 5 and 8 in the empty set. The blue arrows emphasize the invariant P_2 .

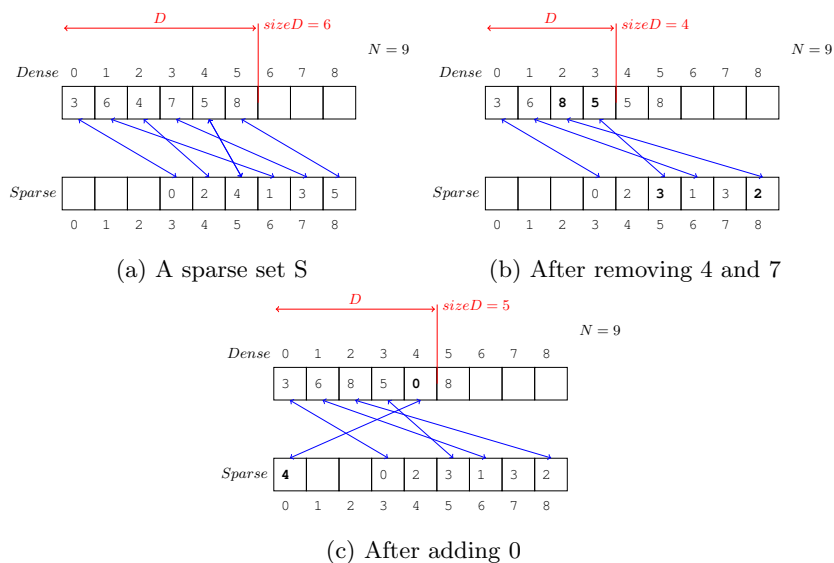


Fig. 1: A sparse set and some operations

Checking if an element v belongs to the sparse set D simply consists in the evaluation of the expression $0 \leq Sparse[v] < sizeD \ \&\& \ Dense[Sparse[v]] = v$. Removing an element consists in replacing v in $Dense$ with the last element e of the $Dense$ effective part ($e = Dense[sizeD - 1]$), decrementing $sizeD$ and updating $Sparse[e]$. This operation is illustrated in Fig. 1b: 4 and 7 are removed in this order from the sparse set represented in Fig. 1a. We can see two occurrences of both 8 and 5 but their presence in $Dense[sizeD..N]$ does not matter.

Inserting an element v is implemented as follows: put v in $Dense$ at the position $sizeD$, update $Sparse[v]$ with $sizeD$ and increment $sizeD$. In Fig. 1c, 0 has been inserted in the sparse set represented in 1b.

Clearing the sparse set, that is making it represent the empty set, is very efficient, just set $sizeD$ to 0. The cardinality of a sparse set is exactly the value of $sizeD$. All the previous operations are constant-time. Operations like forall, exists, union, intersection, equality only require to explore the elements in the effective part of $Dense$, and are thus in $O(sizeD)$.

3 Why3 and WhyML

Why3 [12] is a platform for deductive program verification that provides a specification and programming language called WhyML. It relies on external automated and interactive theorem provers to discharge automatically generated verification conditions (VC). The SMT provers Alt-ergo, CVC4 and Z3 are used here. Transformations, aka tactics, are also provided, making Why3 an interactive proof environment. Why3 supports modular verification and includes some mechanisms for managing modularity, abstraction and genericity [13].

WhyML allows the user to write functional or imperative programs featuring polymorphism, algebraic data types, pattern-matching, exceptions, mutable variables, arrays, etc. These programs can be specified by using contracts (pre- and post- conditions) and assertions (e.g. variants, loop invariants). User-defined types with invariants can be introduced, invariants are verified at the function call boundaries. Furthermore to prevent logical inconsistencies, Why3 generates a verification condition to ensure that such a type is inhabited. To help the verification, a witness can be explicitly given by the user (by clause in Fig. 3). The `old` operator can be used inside post-conditions to refer to the value of a term at the call program point.

Correct-by-construction OCaml (and, more recently, C) programs can be automatically extracted from verified WhyML programs. More detail is provided throughout the paper as necessary.

4 Formal Verification of Sparse Sets

This section deals with the data structure as it is described in Briggs and Torzon's paper [4]. We provide a WhyML specification and an implementation of the data structure and its operations.

4.1 Abstract Specification

We start with a high-level module that contains the abstract specification of type τ and operations on that type, where τ is the type of subsets of an interval of natural numbers (beginning at 0 as in [4]). Fig. 2 contains an excerpt of that module. The type τ here is specified as a record with two fields only used for specification: the size n of the support universe and `setD` which has to be understood as the high-level model of the data structure. The `fset` logical type constructor is defined in the module `set.FsetInt` of the standard library. Set mathematical symbols that appear in the contracts are used here and in the rest of the paper to denote the mathematical set operations acting on mathematical sets also defined in the module `set.FsetInt`. The `==` infix operator is the mathematical set equality. The `writes` clause in a contract indicates that the corresponding function updates its argument. The operations are implemented in the refining module that also provides a full definition for the type. We describe this refining module in the next subsection.

```

module FiniteNatSet
use int.Int
use set.FsetInt

type t = abstract {n : int ;
                   mutable setD : fset int ; }
invariant {setD  $\subseteq$  (interval 0 n)}

val empty_set (nn : int) : t
requires {0<=nn}
ensures {result.setD =  $\emptyset$ }
ensures {result.n = nn}

val member (v : int) (a : t) : bool
requires {0<=v}
ensures {result = v  $\in$  a.setD}

val cardinal_sparse (a : t) : int
ensures {result = |a.setD|}

val add (v : int) (a : t) : unit
requires {0<=v<a.n}
ensures {a.setD == (old a.setD)  $\cup$  {v} }
writes {a.setD}

val remove (v : int) (a : t) : unit
requires {0<=v<a.n}
ensures {a.setD == (old a.setD) - {v}}
writes {a.setD}

end

```

Fig. 2: Abstract specification

4.2 Concrete Implementation

The abstract type t is implemented as the record type `tsparse` (see Fig.3) whose fields are the size of the universe n , the two mutable arrays `dense` and `sparse`, the mutable bound `sizeD` such that the subarray `dense[0..sizeD-1]` contains the current elements of the sparse set and the ghost mathematical and abstract model `setD`. Why3 will generate verification conditions to ensure that the concrete implementation respects the abstract specification.

This record type definition is constrained by invariant properties: the length of both arrays is n which is a positive number, contents are belonging to the integer range $0..n-1$ (Inv1), `sizeD` is between 0 and n (Inv2), the two arrays must be consistent for those elements in the set (Inv3) (P_2 in Sect. 2). Furthermore the last property, Inv5, relates the abstract model with the concrete representation as in the property P_1 of Sect. 2.

In [4], Briggs and Torczon emphasize the fact that the two arrays do not require to be initialized when allocated. In the solution given by Filliâtre and Paskevich to the formal verification of sparse arrays, the arrays are not initialized too [11]. They specify a non initialized memory with the help of a `malloc` function. In our implementation we initialize the arrays `dense` and `sparse` with a negative value (-1) when they are created. We could have reused this approach in our formalization, but we did not in order to stay in line with the variant developed in the next section.

```

let constant initval : int = -1
predicate dom_ran (a : array int) (n: int) =
  0 <= n && a.length = n && forall i. 0<=i<n -> initval<=a[i]< n

type tsparse =      { n : int;
                    mutable dense: array int;
                    mutable sparse: array int;
                    mutable sizeD: int;
                    mutable ghost setD : fset int; }

invariant {
  (*Inv1*)  dom_ran dense n && dom_ran sparse n &&
  (*Inv2*)  0 <= sizeD <= n &&
  (*Inv3*)  (forall i:int. 0 <= i < sizeD ->
            (0 <= dense[i] && sparse[dense[i]]=i)) &&
  (*Inv4*)  setD ⊆ (interval 0 n) &&.
  (*Inv5*)  forall x: int. 0<= x < n -> (x ∈ setD <->
            (0 <= sparse[x] < sizeD && dense[sparse[x]] = x))
}
by {n = 0; dense = make 0 initval; sparse = make 0 initval;
     sizeD = 0; setD = ∅}

```

Fig. 3: WhyML type of a sparse set

The code of the operations on sparse sets are the straightforward translation of the algorithms in [4], except for the supplementary ghost code (e.g. the last statement in `remove_sparse`) which updates the abstract model `a.setD`. The deletion operation, named here `remove_sparse`, is shown in Fig. 4.

In addition to the previous constant-time operations, the following functions have been implemented and verified:

- `forall`: check if all the elements of the sparse set satisfy a predicate (linear with respect to the number of elements in `a.setD`);
- `exists`: check if one element of the sparse set satisfies a predicate (linear with respect to the number of elements in `a.setD`);
- `tolist`: return the list of elements (linear with respect to the number of elements in `a.setD`);
- `copy`: create a copy of a sparse set (linear wrt the number of elements in `a.setD`);
- `union of 2 sparse sets` : create a new sparse set (linear wrt the number of elements in each set);
- `in place union`: update the first argument required to have the largest universe (linear wrt the number of elements in the second argument);
- `intersection of 2 sparse sets`: create a new sparse set (linear wrt the number of the smallest set).

Their deductive verification required to invent and add some formal annotations like loop invariants, ghost code and lemma functions. A typical example is the implementation of `cardinal_sparse` illustrated in Fig. 4. Its code is very simple since the number of elements in the sparse set `a` is exactly `a.sizeD` but a lemma-function, `cardinal_sizeD`, is used to prove the function’s contract as a lemma that will be provided to the provers. The latter states that `cardinal a.setD = a.sizeD` by going through the dense array up to `sizeD` and gathering and counting its elements.

VCS for the functions concern the conformance of the code to the post-condition and also to the invariant attached to the `tsparse` type.

4.3 Proofs

The proof of all the VCs are done automatically using three automatic provers, CVC4, Alt-Ergo and Z3, using the strategy `Auto Level 2`³. Statistics per prover, number of proofs, time (minimum/maximum/average) in seconds, are recorded in Fig. 5.

4.4 Extraction of OCaml Executable Code

To extract OCaml executable code from this development, we modified the previous WhyML code to use machine integers and mathematical integers. In our

³ and only one assertion in the lemma function about the cardinality.

```

let remove_sparse (v : int) (a : tsparse)
requires {0<=v<a.n}
ensures {a.setD == (old a.setD) - {v}}
=
let i = a.sparse[v] in
if 0 <= i < a.sizeD && a.dense[i]=v then
  let e = a.dense[a.sizeD - 1] in
  a.dense[i] <- e ; a.sparse[e] <- i ;
  a.sizeD <- a.sizeD - 1;
  a.setD <- a.setD - {v}

(* a lemma function to help the verification*)
let lemma cardinal_sizeD (a : tsparse)
ensures {|a.setD| = a.sizeD}
=
let ghost ref s = FsetInt.empty in
let ghost ref nb = 0 in
for i = 0 to a.sizeD - 1 do
  invariant {forall x:int. (exists j. 0<=j<i && x = a.dense[j]) <-> x ∈ s}
  invariant {nb = |s| && nb = i}
  s <- s ∪ {a.dense[i]};
  nb <- nb + 1
done ;
assert {a.setD == s && nb = a.sizeD }

let cardinal_sparse (a : tsparse) : int
ensures {result = |a.setD|}
=
return a.sizeD

```

Fig. 4: Implementation of some sparse set operations in WhyML

<i>Prover</i>	<i>nb.proofs</i>	<i>min.time(s)</i>	<i>max.time(s)</i>	<i>av.time(s)</i>
Z3 4.8.9	1	0.05	0.05	0.05
Alt-Ergo 2.5.1	19	0.09	1.15	0.54
CVC4 1.6	205	0.03	8.18	0.18

Fig. 5: Statistics per prover: number of proofs, time (minimum/maximum/average) in seconds

case it requires only syntactical modifications regarding the type of integer variables and arrays and some insertions of coercions between machine integers and mathematical integers in the logical assertions. The proofs remain all automatic.

This data structure is also often proposed as a bounded data structure, in which the set is constrained to have at most a given cardinality m . We can find several implementations of this variant on the Web. In that case the length of the dense array is m . The abstract type \mathfrak{t} and the concrete type $\mathfrak{t}\text{sparse}$ are modified to take into account this maximal capacity. Some functions (e.g. `add` and `union`) are also concerned with this limit. This new requirement does not bring any difficulty for the verification. We have implemented this variant in WhyML and verified it with Why3. When machine integers are used, the `union` function requires an additional pre-condition for not going to an overflow.

5 Backtrackable Sparse Sets as Domains

In this section we focus on a variant of sparse sets used in some constraint solvers (e.g. MiniCP [18], OsCaR [20]) to represent the domain of an integer variable, i.e. the finite set of possible values for that variable [15]. In such a context, to find a solution to a collection of constraints on some variables, or to show that the problem is unsatisfiable, the use case is as follows: for a variable X , initialize $\text{Domain}(X) = 0..N - 1$, for some N , then propagate constraints to prune $\text{Domain}(X)$, then set $\text{Domain}(X)$ to a singleton containing a value of the pruned domain, propagate again, etc., backtrack if necessary. Thus, once the domain is initialized, there is no need to add any value, only deletions are performed. The advantage of sparse sets, as we have seen, is that membership and deletion operations can be performed in constant time. Furthermore, with a simple variation, these data structures are easy to restore when exploring solutions in an imperative setting, making backtracking cheap. Even if they are not used in constraint solving, we keep in our verified implementation the `add` and `union` operations but we will have to take care of the fact that they break reversibility. In the rest of the paper, to refer to this variant, we sometimes use the expression *sparse sets as domains* or shortly *domains*.

5.1 Algorithmic Variant

In this variant, the property P_2 is enforced for every value in Dense (not only in $\text{Dense}[0..\text{sizeD}-1]$): $\text{Sparse}[\text{Dense}[i]] = i$ for all $i \in 0..N - 1$, called now P'_2 . Checking the membership of value v becomes trivial: just check $\text{Sparse}[v] < \text{sizeD}$. Removing an element v now consists of swapping v with the last element in Dense , decrementing sizeD and also updating Sparse . An example is shown in Fig. 6. As pointed out in [15], the values in $\text{dense}[\text{sizeD}..N - 1]$ are not changed by any operation, in particular by a deletion. Let us call this property P_3 . This property can easily be added as an additional post-condition of the `remove` operation. The other operations remain the same (even if `add` and `union` are not used in constraint solving). We introduce a new function `bind`, which

takes an argument v and reduces the set to the singleton $\{v\}$. It is useful in the context of constraint solving, to bind the value of a variable when exploring the search space. Its behaviour is very similar to `remove`: v is swapped with the last element in `sparse`, `dense` is updated accordingly, and `sizeD` is set to 1. Illustrations are given in Fig. 6.

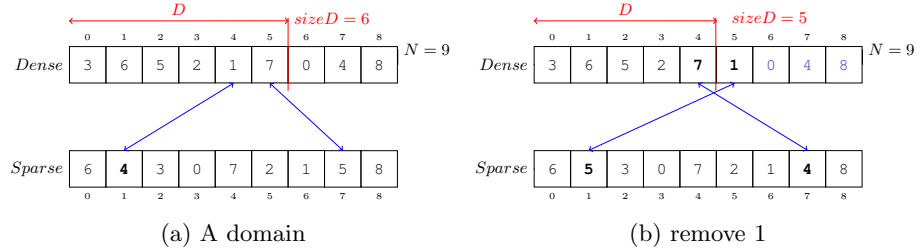


Fig. 6: A sparse set as domain and a deletion

Sparse sets in this variant are now easily backtrackable (or reversible), the only element to be stored and restored being the value of `sizeD`. Fig. 7 illustrates this with a simple example. Let a be the sparse set in Fig. 7a denoting the set D_0 . We store the current value of `a.sizeD`, which is 6. Then we remove 1, 6 and 3 from a , whose resulting value is described on Fig. 7b. To restore the initial situation, it is sufficient to set `a.sizeD` to the value previously stored, i.e. 6, see Fig. 7c. We recover exactly the same mathematical set of elements, even if the value of the two arrays is different from the value in Fig. 7a. The behaviour is the same when backtracking after a `bind` operation, or a destructive intersection operation, or a combination of all of these. However insertion and destructive union operations break this possibility, we keep them but after their use, all checkpoint information is lost. We introduce the operation `undo` to come back to a previously met situation, its algorithmic content is very simple but its specification needs more work.

5.2 WhyML formalization

We follow the same approach with an abstract specification (module `FiniteNat-Set`) and a concrete implementation. To take into account P_2' , we have to change the type invariant of the `tsparse` type. To specify `undo`, the modification is deeper and impacts both the abstract type `t` and the concrete type `tsparse`.

Abstract Specification The type definition of the abstract type `t` is shown on Fig. 8. We introduce an additional abstract variable, `states`, of type `array (option (fset int))`⁴, which stores the different successive states of the

⁴ Since it is only used for specification and proof, it would be better to use a sequence or a function rather than an array, which is usually reserved for code.

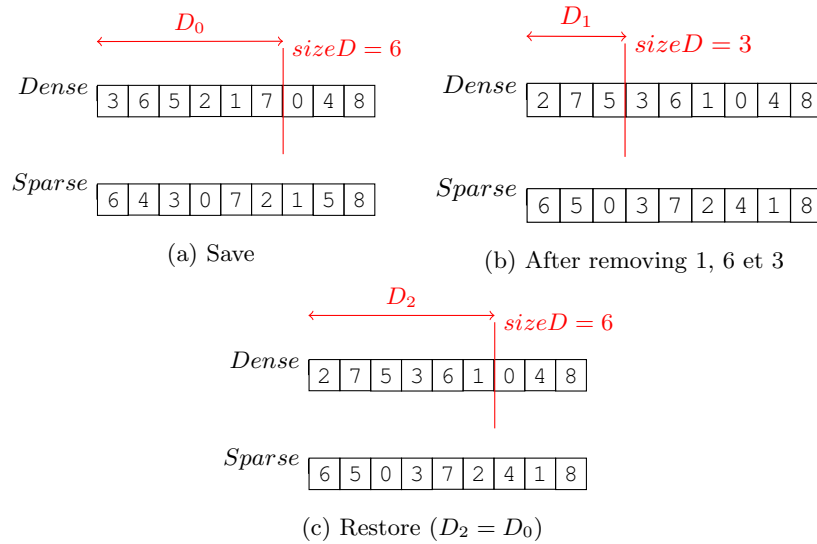


Fig. 7: Backtracking on a sparse set as domain

set, i.e. the successive mathematical models. It is defined as an array of length $n+1$, indexed by all the possible cardinalities of `setD` and containing either `None`, meaning that the corresponding cardinality is not reachable by undoing operations or `Some s` where s is a mathematical set. Each time an operation modifying the sparse set is performed, its current mathematical model is stored in `states[|setD|]`. Because of the `bind` operation, `states` may have *holes* (i.e. `None`) between two registered models. The predicate `valid_states` specifies the well-formedness of the array `states`.

Fig. 8 also contains the abstract specification of the `undo` operation. It takes an argument that is the cardinality of the set to which to come back. This one must be greater than the cardinality of the current set and the corresponding state must have been encountered in the past, so its mathematical model must have been registered in `states`. The post-condition specifies that after calling the operation the model is the one stored in `states[p]` and the previous states have not changed for cardinalities greater than p . The new value of `states` is well-formed, which is implicit since the type invariant associated with `t` must be preserved.

Concrete Implementation The `tsparse` type is adapted in the same way as the abstract type `t` (see Fig. 9). An additional ghost variable `states` is introduced and constrained to be well-formed. Furthermore a new property, `Inv10`, makes the connection between `states` and `sparse`: if s is the mathematical set registered in `states` at index i , then its elements are exactly those that are in `dense[0,i]`. Again as a property in the type invariant, it must be preserved

```

predicate valid_states (states : array (option (fset int)))
    (setD : fset int) (n : int) =
  states.length = n+1 &&
  states[|setD|] = Some setD &&
  (forall i. 0 <= i < |setD| -> states[i] = None) &&
  (forall i:int, s : fset int. 0<=i<=n -> states[i] = Some s ->
    s ⊆ (interval 0 n) && |setD| <= i && setD ⊆ s)

type t = abstract {n : int ; mutable setD : fset int ;
    mutable states: array (option (fset int));
    }
invariant {subset setD (interval 0 n) &&
  valid_states states setD n
  }

val remove (v : int) (a : t) : unit
requires {0<=v<a.n && v ∈ a.setD}
ensures {a.setD == remove v (old a.setD)}
ensures {forall i. 0 <= i <= a.n -> i ≠ |a.setD| -> a.states[i] = (old a).states[i]}

val add (v : int) (a : t) : unit
requires {0<=v<a.n && v ∉ a.setD}
ensures {a.setD == add v (old a.setD)}
ensures {forall i. 0 <= i <= a.n -> i ≠ |a.setD| -> a.states[i] = None}

val undo (a : t) (p : int) : unit
requires {cardinal (a.setD) < p <= a.n}
requires {exists s. a.states[p] = Some s}
ensures { a.setD == from_option ((old a).states[p])}
ensures {forall i. p < i <= a.n -> a.states[i] = (old a).states[i]}

```

Fig. 8: Abstract type t and undo abstract specification

by any operation modifying an argument of type `tsparse`. Furthermore *Inv3* is modified to take into account that `dense` and `sparse` are now inverse and *Inv5* is simplified.

The code of the undo operation is shown in Fig. 9. Its contract is similar to that of the abstract specification. Its computational part is only the last statement, the rest is some ghost code to update the model (`a.setD` and `a.states`). In particular, to maintain the invariant, all states between `a.sizeD` and `p` are deleted in `a.states`, thanks to the `fill` operation. The operations for removing and inserting an element are also illustrated in that figure. The computational part is composed of the two first statements. Besides the modification of the `a.setD` abstract set, the ghost code updates the `a.states` array: the former operation just stores the current state while the latter also erases all the previous stored models.

Proofs The proof of all the VCs are done automatically using two automatic provers, CVC4 and Alt-Ergo using the strategy Auto Level 2. Statistics per prover, number of proofs, time (minimum/maximum/average) in seconds, are recorded in Fig. 10.

6 Some Experimentations

OCaml code was extracted from the WhyML models (using machine integers) of all the variants we have developed. We implemented a naive implementation of the Erathosthenes Sieve algorithm following a Web article⁵ presenting sparse set implementations in C++, using three of our variants of sparse sets: sparse sets *à la* Briggs and Torczon with and without limited capacity and sparse sets as domains. We also implemented this algorithm using the OCaml standard library module `Set` and an implementation of sets as hash tables⁶.

This algorithm performs many insertions, deletions, membership tests and a final call to the operation that computes the cardinality of the sparse set that, at the end, contains the prime numbers up to P , the parameter of the algorithm. In our experimentation whose results are shown on Fig. 11, P varies from one hundred to one million. On the x-axis are the execution times in seconds and on the y-axis the values of P .

On this example sparse sets in their three versions outperform the two set other representations but it is a bit unfair since we are comparing a mutable representation with functional ones. Regarding the three variants, they are equivalent, maintaining the links for removed elements do not impact significantly the execution time.

⁵ <https://www.codeproject.com/Articles/859324/Fast-Implementations-of-Sparse-Sets-in-Cplusplus>

⁶ <https://github.com/backtracking/hashset>

```

type tsparse =      { n : int;
                      mutable dense: array int;
                      mutable sparse: array int;
                      mutable sized: int;
                      mutable ghost setD: fset int;
                      mutable ghost states: array (option (fset int));
                      }

invariant {
  (*Inv1 *)  dom_ran dense n && dom_ran sparse n &&
  (*Inv2 *)  0 <= sized <= n &&
  (*Inv3' *) (forall i:int. 0 <= i < sized && 0<=v<n->
             (dense[i]=v <-> sparse[v]=i))
  (*Inv4 *)  setD ⊆ (interval 0 n) &&.
  (*Inv5' *) (forall x: int. 0<= x < n ->
             (x ∈ setD <-> sparse[x] < sized)) &&
  (*Inv6 *)  states.length = n+1 &&
  (*Inv7 *)  states[sized] = Some setD &&
  (*Inv8 *)  (forall i. 0 <= i < sized -> states[i] = None) &&
  (*Inv9 *)  (forall i, s. 0<=i<=n -> states[i] = Some s ->
             (s ⊆ (interval 0 n) && sized <= i && setD ⊆ s)) &&
  (*Inv10 *) (forall i, s. 0<=i<=n -> states[i] = Some s ->
             (forall x. 0<=x<n -> (sparse[x]<i <-> mem x s)))

let undo_sparse (a : tsparse) (p : int) : unit
...
=
let ghost v = a.states[p] in
  a.setD <- from_option v ;
  fill a.states 0 p None;
  a.sized <- p

let remove_sparse (v : int) (a : tsparse)
...
=
swap_two_arrays a.dense a.sparse a.n a.sparse[v] (a.sized - 1);
a.sized <- a.sized - 1;
a.setD <- remove v a.setD;
a.states[a.sized] <- Some a.setD

let add_sparse (v : int) (a : tsparse)
...
=
swap_two_arrays a.dense a.sparse a.n a.sparse[v] a.sized;
a.sized <- a.sized + 1;
a.setD <- add v a.setD;
fill a.states 0 (a.n + 1) None ;
a.states[a.sized] <- Some a.setD;

```

Fig. 9: Concrete Implementation of Domains

<i>Prover</i>	<i>nb.proofs</i>	<i>min.time(s)</i>	<i>max.time(s)</i>	<i>av.time(s)</i>
Alt-Ergo 2.5.1	37	0.01	1.09	0.29
CVC4 1.6	148	0.05	0.93	0.12

Fig. 10: Statistics per prover: number of proofs, time (minimum/maximum/average) in seconds

Our extracted code of sparse sets as domains has been used with a simple sudoku solver originally written in OCaml by Filiâtre⁷. That example intensively uses backtracking and thus the undo operation. It has been evaluated on a large number of sudoku puzzles.

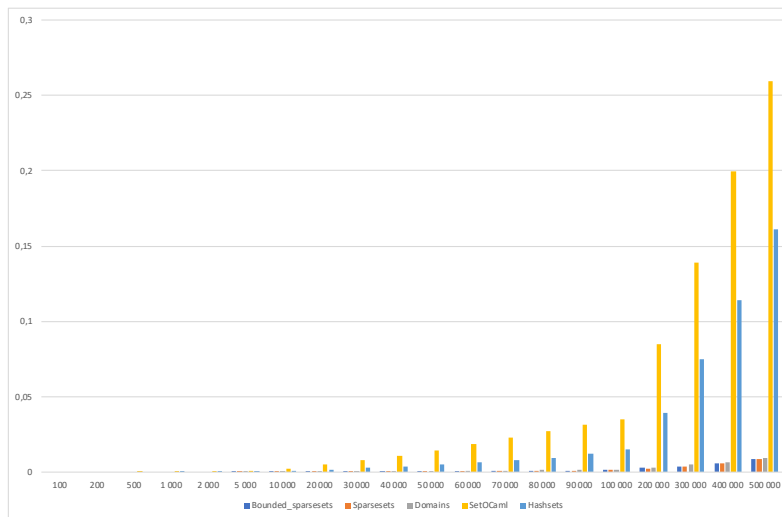


Fig. 11: Comparison of execution times on a naive implementation of Erathos-thenes Sieve algorithm

7 Conclusion

In this paper we presented the formal Why3 development for sparse sets and for sparse sets as domains used in constraint solvers. The former refines and

⁷ <https://github.com/backtracking/ocaml-bazaar/blob/main/sudoku.ml>

extends a partial solution for a more general data structure, sparse arrays, done by Filliâtre and Paskevich some years ago. The latter is a variant of the former, but as far as we know it is the first formalization of this backtrackable data structure that allows the representation of domains of integer variables. We have extracted efficient OCaml code from these formally verified models, which we have experimented on simple test cases and the Erathosthenes Sieve algorithm. One perspective of this work is the extraction of C code.

The technique used to be able to specify and prove the undo operation has implications for the whole formal development. It allows to use WhyML and the deductive verification engine of Why3 to prove a property that involves more than a pre- and a post- state, and is close to a dynamic or temporal property.

In the case of very sparse sets or domains, using an array to implement the sparse structure is not optimal in terms of memory space. The data structure could be made more interesting by using another fast access structure, e.g. a hashmap (idea also suggested in [15]). So we could also suggest extending our current work to use such an alternative. It would be more interesting to make this sparse structure a generic parameter of the formalization in order to choose the right implementation *à la carte*.

As future work, we would also like to integrate sparse sets as domains in a finite domain constraint solver, e.g. in CoqBinFD, a formally verified constraint solver formally verified in Coq [6] or in FaCile, an OCaml constraint library [5].

Acknowledgements The author would like to thank the Why3 development team at LMF and also Jean-Paul Bodeveix, for their interest and helpful comments during an initial presentation of this work. She also thanks the anonymous reviewers for their suggestions and careful reading.

References

1. J.-R. Abrial. *The B-Book, Assigning Programs to Meanings*. Cambridge University Press, 1996.
2. J.-R. Abrial. *Modeling in Event-B - System and Software Engineering*. Cambridge University Press, 2010.
3. A. V. Aho and J. E. Hopcroft. *The Design and Analysis of Computer Algorithms*. Addison-Wesley Longman Publishing Co., Inc., USA, 1st edition, 1974.
4. P. Briggs and L. Torczon. An efficient representation for sparse sets. *LOPLAS*, 2(1-4):59–69, 1993.
5. P. Brisset and N. Barnier. FaCiLe : a Functional Constraint Library. In *CICLOPS 2001, Colloquium on Implementation of Constraint and LOGic Programming Systems*, Paphos, Cyprus, 2001.
6. M. Carlier, C. Dubois, and A. Gotlieb. A certified constraint solver over finite domains. In *Formal Methods (FM 2012)*, volume 7436 of *LNCS*, pages 116–131, Paris, France, 2012.
7. M. Cristiá and C. Dubois. Comparing EventB, {log} and Why3 Models of Sparse Sets. In *35es Journées Francophones des Langages Applicatifs (JFLA 2024)*, Saint-Jacut-de-la-Mer, France, Jan. 2024.
8. M. Cristiá and G. Rossi. {log}: set formulas as programs. *Rend. Ist. Mat. Univ. Trieste*, 53:24, 2021. Id/No 23.
9. J. Demeulenaere, R. Hartert, C. Lecoutre, G. Perez, L. Perron, J. Régis, and P. Schaus. Compact-table: Efficiently filtering table constraints with reversible sparse bit-sets. In M. Rueher, editor, *Principles and Practice of Constraint Programming - 22nd International Conference, CP 2016, Toulouse, France, September 5-9, 2016, Proceedings*, volume 9892 of *Lecture Notes in Computer Science*, pages 207–223. Springer, 2016.
10. J.-C. Filliâtre and P. Letouzey. Functors for proofs and programs. In D. A. Schmidt, editor, *Programming Languages and Systems, 13th European Symposium on Programming, ESOP 2004, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2004, Barcelona, Spain, March 29 - April 2, 2004, Proceedings*, volume 2986 of *Lecture Notes in Computer Science*, pages 370–384. Springer, 2004.
11. J.-C. Filliâtre and A. Paskevich. Why3 version of the sparse arrays example, the first example of the vacid-0 benchmarks, https://toccata.gitlabpages.inria.fr/toccata/gallery/vacid_0_sparse_array.en.html.
12. J.-C. Filliâtre and A. Paskevich. Why3 - where programs meet provers. In M. Felleisen and P. Gardner, editors, *22nd European Symposium on Programming, ESOP 2013, Held as Part of ETAPS 2013, Rome, Italy, Proceedings*, volume 7792 of *LNCS*, pages 125–128. Springer, 2013.
13. J.-C. Filliâtre and A. Paskevich. Abstraction and genericity in why3. In T. Margaria and B. Steffen, editors, *Leveraging Applications of Formal Methods, Verification and Validation: Verification Principles - 9th International Symposium on Leveraging Applications of Formal Methods, ISOFA 2020, Rhodes, Greece, October 20-30, 2020, Proceedings, Part I*, volume 12476 of *Lecture Notes in Computer Science*, pages 122–142. Springer, 2020.
14. L. Lamport. *Specifying Systems, The TLA+ Language and Tools for Hardware and Software Engineers*. Addison-Wesley, 2002.
15. V. Le Clément de Saint-Marcq, P. Schaus, C. Solnon, and C. Lecoutre. Sparse-Sets for Domain Implementation. In *CP workshop on Techniques for Implementing Constraint programming Systems (TRICS)*, pages 1–10, 2013.

16. A. Ledein and C. Dubois. Facile en coq : vérification formelle des listes d'intervalles. In *31ème Journées Francophones des Langages Applicatifs*, 2019.
17. K. R. M. Leino and M. Moskal. Vacid-0: Verification of ample correctness of invariants of data-structures, edition 0. In *Proceedings of Tools and Experiments Workshop at VSTTE*, 2010.
18. L. Michel, P. Schaus, and P. Van Hentenryck. Minicp: a lightweight solver for constraint programming. *Mathematical Programming Computation*, 13(1):133–184, 2021.
19. T. Nipkow, J. Blanchette, M. Eberl, A. Gómez-Londoño, P. Lammich, C. Sternagel, S. Wimmer, and B. Zhan. *Functional Algorithms, Verified!* Freely downloadable <https://functional-algorithms-verified.org>, 2023.
20. P. Schaus and R. D. Landtsheer. Oscar user-guide, 2016. Available from <http://oscarlib.org>.
21. J. T. Schwartz, R. B. K. Dewar, E. Dubinsky, and E. Schonberg. *Programming with Sets - An Introduction to SETL*. Texts and Monographs in Computer Science. Springer, 1986.