

Chapter 2

Lock-free dynamic hash tables with open addressing

This chapter has been published as [21]. Its extended abstract appears in [20].

2.1 Introduction

We are interested in efficient, reliable, parallel algorithms. The classical synchronization paradigm based on mutual exclusion is not most suited for this, since mutual exclusion often turns out to be a performance bottleneck, and failure of a single process can force all other processes to come to a halt. This is the reason to investigate lock-free or wait-free concurrent objects, see e.g. [6, 27, 28, 30, 38, 46, 48, 55, 64, 66, 67, 69].

Lock-free and wait-free objects

A *concurrent object* is an abstract data type that permits concurrent operations that appear to be atomic [27, 52, 69]. The easiest way to implement concurrent objects is by means of mutual exclusion, but this leads to blocking when the process that holds exclusive access to the object is delayed or stops functioning.

The object is said to be *lock-free* if any process can be delayed at any point without forcing any other process to block and when, moreover, it is guaranteed that always some process will complete its operation in a finite number of steps, regardless of the execution speeds of the processes [6, 28, 48, 64, 69]. The object is said to be *wait-free* when it is guaranteed that any process can complete any operation in a finite number of steps, regardless of the speeds of the other processes [27].

We regard “non-blocking” as synonymous to “lock-free”. In several recent papers, e.g. [67], the term “non-blocking” is used for the first conjunct in the above definition of lock-free. Note that this weaker concept does not in itself guarantee progress. Indeed, without real blocking, processes might delay each other arbitrarily without getting closer to completion of their respective operations. The older literature [2, 6, 30] seems to suggest that originally “non-blocking” was used for the stronger concept, and lock-free for the weaker one. Be this as it may, we use lock-free for the stronger concept.

Concurrent hash tables

The data type of hash tables is very commonly used to efficiently store huge but sparsely filled tables. Before 2003, as far as we know, no lock-free algorithm for hash tables had been proposed. There were general algorithms for arbitrary wait-free objects [3, 6, 7, 27, 35, 36], but these are not very efficient. Furthermore, there are lock-free algorithms for different domains, such as linked lists [69], queues [68] and memory management [30, 38].

In this chapter we present a lock-free algorithm for hash tables with open addressing that is in several aspects wait-free. The central idea is that every process holds a pointer to a hash table, which is the current one if the process is not delayed. When the current hash table is full, a new hash table is allocated and all active processes join in the activity to transfer the contents of the current table to the new one. The consensus problem of the choice of a new table is solved by means of a test-and-set register. When all processes have left the obsolete table, it is deallocated by the last one leaving. This is done by means of a compare-and-swap register. Measures have been taken to guarantee that actions of delayed processes are never harmful. For this purpose we use counters that can be incremented and decremented atomically.

After the initial design, it took us several years to establish the safety properties of the algorithm. We did this by means of the proof assistant PVS [63]. Upon completion of this proof, we learned that a lock-free resizable hash table based on chaining was proposed in [66]. We come back to this below.

Our algorithm is lock-free and some of the subtasks are wait-free. We allow fully parallel *insertion*, *assignment*, *deletion*, and *finding* of elements. Finding is wait-free, the other three are not. The primary cause is that the process executing it may repeatedly have to execute or join a migration of the hash table. Assignment and deletion are also not wait-free when other processes repeatedly assign to the same address successfully.

Migration is called for when the current hash table is almost filled. This occurs when the table has to grow beyond its current upper bound, but also for maintenance after many insertions and deletions. The migration itself is wait-free, but, in principle, it is possible that a slow process is unable to access and use a current hash table since the current hash table is repeatedly replaced by faster processes.

Migration requires subtle provisions, which can be best understood by considering the following scenario. Suppose that process A is about to (slowly) insert an element in a hash table H_1 . Before this happens, however, a fast process B has performed migration by making a new hash table H_2 , and copying the content from H_1 to H_2 . If (and only if) process B did not copy the insertion of A , A must be informed to move to the new hash table, and carry out the insertion there. Suppose a process C comes into play also copying the content from H_1 to H_2 . This must be possible, since otherwise B can stop copying, blocking all operations of other processes on the hash table, and thus violating the lock-free nature of the algorithm. Now the value inserted by A can but need not be copied by both

B and/or C . This can be made more complex by a process D that attempts to replace H_2 by H_3 . Still, the value inserted by A should show up exactly once in the hash table, and it is clear that processes should carefully keep each other informed about their activities on the tables.

Performance, comparison, and correctness

Actually, only an extra check is required in the main loop of the main functions, one extra bit needs to be set when writing data in the hashtables and at some places a write operation has been replaced by a compare-and-swap, which is more expensive. For ordinary operations on the hashtable, this is the only overhead and therefore a linear speed up can be expected on multiprocessor systems. The only place where no linear speed up can be achieved is when copying the hashtable. Especially, when processes have widely different speeds, a logarithmic factor may come into play (see algorithms for the write all problem [24, 46]). Indeed, initial experiments indicate that our algorithm is as efficient as sequential hash tables. It seems to require on average only constant time for insertion, deletion or accessing of elements.

Some differences between our algorithm and the algorithm of [66] are clear. No formally verified correctness proof was given for the algorithm in [66]. In our algorithm, the hashed values need not be stored in dynamic nodes if the address-value pairs (plus one additional bit) fit into one word. Our hash table can shrink whereas the table of bucket headers in [66] cannot shrink. A disadvantage of our algorithm, due to its open addressing, is that migration is needed as maintenance after many insertions and deletions.

An apparent weakness of our algorithm is the worst-case space complexity in the order of $\mathcal{O}(PM)$ where P is the number of processes and M is the size of the table. This only occurs when many of the processes fail or fall asleep while using the hash table. Failure while using the hash table can be made less probable by adequate use of the procedure “releaseAccess”. This gives a trade-off between space and time since it introduces the need of a corresponding call of “getAccess”. When all processes make ordinary progress and the hash table is not too small, the actual memory requirement is $\mathcal{O}(M)$.

The migration activity requires worst-case $\mathcal{O}(M^2)$ time for each participating process. This only occurs when the migrating processes tend to choose the same value to migrate and the number of collisions is $\mathcal{O}(M)$ due to a bad hash function. This is costly, but even this is in agreement with wait-freedom. The expected amount of work for migration for all

processes together is $\mathcal{O}(M)$ when collisions are sparse, as should be the case when migrating to a hash table that is sufficiently large.

A true problem of lock-free algorithms is that they are hard to design correctly, which even holds for apparently straightforward algorithms. Whereas human imagination generally suffices to deal with all possibilities of sequential processes or synchronized parallel processes, this appears impossible (at least to us) for lock-free algorithms. The only technique that we see fit for any but the simplest lock-free algorithms is to prove the correctness of the algorithm very precisely, and to verify this using a proof checker or theorem prover.

As a correctness notion, we take that the operations behave the same as for ‘ordinary’ hash tables, under some arbitrary linearization [31] of these operations. So, if a *find* is carried out strictly after an *insert*, the inserted element is found. If *insert* and *find* are carried out at the same time, it may be that *find* takes place before *insertion*, and it is not determined whether an element will be returned.

Our algorithm contains 81 atomic statements. The structure of our algorithm and its correctness properties, as well as the complexity of reasoning about them, makes neither automatic nor manual verification feasible. We have therefore chosen the higher-order interactive theorem prover PVS [11, 63] for mechanical support. PVS has a convenient specification language and contains a proof checker which allows users to construct proofs interactively, to automatically execute trivial proofs, and to check these proofs mechanically.

Overview of this chapter

Section 2.2 contains the description of the hash table interface offered to the users. The algorithm is presented in Section 2.3. Section 2.4 contains a description of the proof of the safety properties of the algorithm: functional correctness, atomicity, and absence of memory loss. This proof is based on a list of around 200 invariants, presented in Appendix A.1, while the relationships between the invariants are given by a dependency graph in Appendix A.2. Progress of the algorithm is proved informally in Section 2.5. Conclusions are drawn in Section 2.6.

2.2 The interface

The aim is to construct a hash table that can be accessed simultaneously by different processes in such a way that no process can passively block another process’ access to the

table.

A hash table is an implementation of (partial) functions between two domains, here called *Address* and *Value*. The hash table thus implements a modifiable shared variable $X : \text{Address} \rightarrow \text{Value}$. The domains *Address* and *Value* both contain special default elements $0 \in \text{Address}$ and $\mathbf{null} \in \text{Value}$. An equality $X(a) = \mathbf{null}$ means that no value is currently associated with the address a . In particular, since we never store a value for the address 0, we impose the invariant

$$X(0) = \mathbf{null} .$$

We use open addressing to keep all elements within the table. For the implementation of the hash table we require that from every value the address it corresponds to is derivable. We therefore assume that some function $ADR : \text{Value} \rightarrow \text{Address}$ is given with the property:

$$\text{Ax1: } v = \mathbf{null} \equiv ADR(v) = 0.$$

Indeed, we need \mathbf{null} as the value corresponding to the undefined addresses and use address 0 as the (only) address associated with the value \mathbf{null} . We thus require the hash table to satisfy the invariant

$$X(a) \neq \mathbf{null} \Rightarrow ADR(X(a)) = a .$$

Note that the existence of ADR is not a real restriction since one can choose to store the pair (a, v) instead of v . When a can be derived from v , it is preferable to store v , since that saves memory.

There are four principle operations: *find*, *delete*, *insert* and *assign*. The first one is to *find* the value currently associated with a given address. This operation yields \mathbf{null} if the address has no associated value. The second operation is to *delete* the value currently associated with a given address. It fails if the address was empty, i.e. $X(a) = \mathbf{null}$. The third operation is to *insert* a new value for a given address, provided the address was empty. So, note that at least one out of two consecutive *inserts* for address a must fail, except when there is a *delete* for address a in between them. The operation *assign* does the same as *insert*, except that it rewrites the value even if the associated address is not empty. Moreover, *assign* never fails.

We assume that there is a bounded number of processes that may need to interact with the hash table. Each process is characterized by the sequence of operations

$$(\text{getAccess} ; (\text{find} + \text{delete} + \text{insert} + \text{assign})^* ; \text{releaseAccess})^\omega$$

A process that needs to access the table, first calls the procedure *getAccess* to get the current hash table pointer. It may then invoke the procedures *find*, *delete*, *insert*, and *assign* repeatedly, in an arbitrary, serial manner. A process that has access to the table can call *releaseAccess* to log out. The processes may call these procedures concurrently. The only restriction is that every process can do at most one invocation at a time.

The basic correctness conditions for concurrent systems are functional correctness and atomicity, say in the sense of [52], chapter 13. Functional correctness is expressed by prescribing how the procedures *find*, *insert*, *delete*, *assign* affect the value of the abstract mapping X in relation to the return value. Atomicity means that the effect on X and the return value takes place atomically at some time between the invocation of the routine and its response. Each of these procedures has the precondition that the calling process has access to the table. In this specification, we use auxiliary private variables declared locally in the usual way. We give them the suffix *S* to indicate that the routines below are the specifications of the procedures. We use angular brackets \langle and \rangle to indicate atomic execution of the enclosed command.

```

proc findS(a : Address \ {0}) : Value =
  local rS : Value;
(fS)   $\langle$  rS := X(a)  $\rangle$ ;
return rS.

```

```

proc deleteS(a : Address \ {0}) : Bool =
  local sucS : Bool;
(dS)   $\langle$  sucS := (X(a) ≠ null) ;
      if sucS then X(a) := null end  $\rangle$  ;
return sucS.

```

```

proc insertS(v : Value \ {null}) : Bool =
  local sucS : Bool ; a : Address := ADR(v) ;
(iS)   $\langle$  sucS := (X(a) = null) ;
      if sucS then X(a) := v end  $\rangle$  ;
return sucS.

```

```

proc assignS( $v : \text{Value} \setminus \{\mathbf{null}\}$ ) =
    local  $a : \text{Address} := \text{ADR}(v)$  ;
(aS)    $\langle \mathbf{X}(a) := v \rangle$  ;
end.

```

Note that, in all cases, we require that the body of the procedure is executed atomically at some moment between the beginning and the end of the call, but that this moment need not coincide with the beginning or end of the call. This is the reason that we do not (e.g.) specify *find* by the single line **return** $\mathbf{X}(a)$.

Due to the parallel nature of our system we cannot use pre- and postconditions to specify it. For example, it may happen that *insert*(v) returns *true* while $\mathbf{X}(\text{ADR}(v)) \neq v$ since another process deletes $\text{ADR}(v)$ between the execution of (iS) and the response of *insert*.

In Section 2.3.4, we provide implementations for the operations *find*, *delete*, *insert*, *assign*. We prove partial correctness of the implementations by extending them with the auxiliary variables and commands used in the specification. So, we regard \mathbf{X} as a shared auxiliary variable and rS and $sucS$ as private auxiliary variables; we augment the implementations of *find*, *delete*, *insert*, *assign* with the atomic commands (fS), (dS), (iS), (aS), respectively. We prove that each of the four implementations executes its specification command always exactly once and that the resulting value r or suc of the implementation equals the resulting value rS or $sucS$ in the specification. It follows that, by removing the implementation variables from the combined program, we obtain the specification. This removal may eliminate many atomic steps of the implementation. This is analogous to removal of stutterings in TLA [49] or abstraction from τ steps in process algebras.

2.3 The algorithm

An implementation consists of P processes along with a set of variables, for $P \geq 1$. Each process, numbered from 1 up to P , is a sequential program comprised of atomic statements. Actions on private variables can be added to an atomic statement, but all actions on shared variables must be separated into atomic accesses. Since auxiliary variables are only used to facilitate the proof of correctness, they can be assumed to be touched instantaneously without violation of the atomicity restriction.

2.3.1 Hashing

We implement function \mathbf{X} via hashing with open addressing, cf. [47, 70]. We do not use direct chaining, where colliding entries are stored in a secondary list, as is done in [66]. A disadvantage of open addressing with deletion of elements is that the contents of the hash table must regularly be refreshed by copying the non-deleted elements to a new hash table. As we wanted to be able to resize the hash tables anyhow, we consider this less of a burden.

In principle, hashing is a way to store address-value pairs in an array (hash table) with a length much smaller than the number of potential addresses. The indices of the array are determined by a hash function. In case the hash function maps two addresses to the same index in the array there must be some method to determine an alternative index. The question how to choose a good hash function and how to find alternative locations in the case of open addressing is treated extensively elsewhere, e.g. [47].

For our purposes it is convenient to combine these two roles in one abstract function *key* given by:

$$\text{key}(a : \text{Address}, l : \text{Nat}, n : \text{Nat}) : \text{Nat} ,$$

where l is the length of the array (hash table), that satisfies

$$\text{Ax2: } 0 \leq \text{key}(a, l, n) < l$$

for all a , l , and n . The number n serves to obtain alternative locations in case of collisions: when there is a collision, we re-hash until an empty “slot” (i.e. **null**) or the same address in the table is found. The approach with a third argument n is unusual but very general. It is more usual to have a function *Key* dependent on a and l , and use a second function *Inc*, which may depend on a and l , to use in case of collisions. Then our function *key* is obtained recursively by

$$\text{key}(a, l, 0) = \text{Key}(a, l) \text{ and } \text{key}(a, l, n + 1) = \text{Inc}(a, l, \text{key}(a, l, n)) .$$

We require that, for any address a and any number l , the first l keys are all different, as expressed in

$$\text{Ax3: } 0 \leq k < m < l \Rightarrow \text{key}(a, l, k) \neq \text{key}(a, l, m) .$$

2.3.2 Tagging of values

As is well known [47], hashing with open addressing needs a special value $\mathbf{del} \in \mathit{Value}$ to replace deleted values.

When the current hash table becomes full, the processes need to reach consensus to allocate a new hash table of new size to replace the current one. Then all values except \mathbf{null} and \mathbf{del} must be migrated to the new hash table. A value that is being migrated cannot be simply removed, since the migrating process may stop functioning during the migration. Therefore, a value being copied must be tagged in such a way that it is still recognizable. This is done by the function old . We thus introduce an extended domain of values to be called $E\mathit{Value}$, which is defined as follows:

$$E\mathit{Value} = \{\mathbf{del}\} \cup \mathit{Value} \cup \{old(v) \mid v \in \mathit{Value}\}.$$

We furthermore assume the existence of functions $val : E\mathit{Value} \rightarrow \mathit{Value}$ and $oldp : E\mathit{Value} \rightarrow \mathit{Bool}$ that satisfy, for all $v \in \mathit{Value}$:

$$\begin{aligned} val(v) &= v & oldp(v) &= false \\ val(\mathbf{del}) &= \mathbf{null} & oldp(\mathbf{del}) &= false \\ val(old(v)) &= v & oldp(old(v)) &= true \end{aligned}$$

Note that the old tag can easily be implemented by designating one special bit in the representation of Value . In the sequel we write \mathbf{done} for $old(\mathbf{null})$. Moreover, we extend the function ADR to domain $E\mathit{Value}$ by $ADR(v) = ADR(val(v))$.

2.3.3 Data structure

A *Hash table* is either \perp , indicating the absence of a hash table, or it has the following structure:

```
size, bound, occ, dels : Nat;
table : array 0 .. size-1 of EValue.
```

The field `size` indicates the size of the hash table, `bound` the maximal number of places that can be occupied before refreshing the table. Both are set when creating the table and remain constant. The variable `occ` gives the number of occupied positions in the table, while the variable `dels` gives the number of deleted positions. If h is a pointer to a hash

table, we write `h.size`, `h.occ`, `h.dels` and `h.bound` to access these fields of the hash table. We write `h.table[i]` to access the i^{th} *EValue* in the table.

Apart from the *current* hash table, which is the main representative of the variable **X**, we have to deal with *old* hash tables, which were in use before the current one, and *new* hash tables, which can be created after the current one.

We now introduce data structures that are used by the processes to find and operate on the hash table and allow to delete hash tables that are not used anymore. The basic idea is to count the number of processes that are using a hash table, by means of a counter **busy**. The hash table can be thrown away when **busy** is set to 0. An important observation is that **busy** cannot be stored as part of the hash table, in the same way as the variables **size**, **occ** and **bound** above. The reason for this is that a process can attempt to access the current hash table by increasing its **busy** counter. However, just before it wants to write the new value for **busy** it falls asleep. When the process wakes up the hash table might have been deleted and the process would be writing at a random place in memory.

This forces us to use separate arrays **H** and **busy** to store the pointers to hash tables and the **busy** counters. There can be $2P$ hash tables around, because each process can simultaneously be accessing one hash table and attempting to create a second one. The arrays below are shared variables.

```

H : array 1 .. 2P of pointer to Hashtable ;
busy : array 1 .. 2P of Nat ;
prot : array 1 .. 2P of Nat ;
next : array 1 .. 2P of 0 .. 2P .

```

As indicated, we also need arrays **prot** and **next**. The variable **next**[i] points to the next hash table to which the contents of hash table **H**[i] is being copied. If **next**[i] equals 0, this means that there is no next hash table. The variable **prot**[i] is used to guard the variables **busy**[i], **next**[i] and **H**[i] against being reused for a new table, before all processes have discarded them.

We use a shared variable **currInd** to hold the index of the currently valid hash table:

```

currInd : 1 .. 2P .

```

Note however that after a process copies **currInd** to its local memory, other processes may create a new hash table and change **currInd** to point to that one.

It is assumed that initially $H[1]$ is pointing to some hash table. The other initial values of the shared variables are given by

$$\begin{aligned} \text{currInd} &= \text{busy}[1] = \text{prot}[1] = 1, \\ H[i] = \text{busy}[i] = \text{prot}[i] &= 0 \text{ for all } i \neq 1, \\ \text{next}[i] &= 0 \text{ for all } i. \end{aligned}$$

2.3.4 Primary procedures

We first provide the code for the primary procedures, which match directly with the procedures in the interface. Every process has a private variable

$$\textit{index} : 1 \dots 2P;$$

containing what it regards as the currently active hash table. At entry of each primary procedure, it must be the case that the variable $H[\textit{index}]$ contains valid information. In section 2.3.5, we provide the procedure *getAccess* with the main purpose to guarantee this property. When *getAccess* has been called, the system is obliged to keep the hash table at *index* stored in memory, even if there are no accesses to the hash table using any of the primary procedures. A procedure *releaseAccess* is provided to release resources, and it should be called whenever the process will not access the hash table for some time.

Syntax

We use a syntax analogous to Modula-3 [25]. We use $:=$ for the assignment. We use the C-operations $++$ and $--$ for atomic increments and decrements. The semicolon is a separator, not a terminator. The basic control mechanisms are

loop .. end is an infinite loop, terminated by **exit** or **return**.
while .. do .. end and **repeat .. until ..** are ordinary loops.
if .. then .. {elsif ..} [else ..] end is the conditional statement.
case .. end is a case statement.

Types are slanted and start with a capital. Shared variables and shared data elements are in typewriter font. Private variables are slanted or in math italic.

The main loop

We model the clients of the hash table in the following loop. This is not an essential part of the algorithm, but it is needed in the PVS description, and therefore provided here.

```

loop
0:   getAccess() ;
     loop
1:   choose call; case call of
       (f, a) with  $a \neq 0 \rightarrow find(a)$ 
       (d, a) with  $a \neq 0 \rightarrow delete(a)$ 
       (i, v) with  $v \neq \mathbf{null} \rightarrow insert(v)$ 
       (a, v) with  $v \neq \mathbf{null} \rightarrow assign(v)$ 
       (r)  $\rightarrow releaseAccess(index); \mathbf{exit}$ 
     end
     end
end

```

The main loop shows that each process repeatedly invokes its four principle operations with correct arguments in an arbitrary, serial manner. Procedure *getAccess* has to provide the client with a protected value for *index*. Procedure *releaseAccess* releases this value and its protection. Note that **exit** means a jump out of the inner loop.

Procedure *find*

Finding an address in a hash table with open addressing requires a linear search over the possible hash keys until the address or an empty slot is found. The kernel of procedure *find* is therefore:

```

n := 0 ;
repeat  $r := h.table[key(a, l, n)] ; n++ ;$ 
until  $r = \mathbf{null} \vee a = ADR(r) ;$ 

```

The main complication is that, when the process encounters an entry **done** (i.e. *old(null)*), it has to join the migration activity by calling *refresh*.

Apart from a number of special commands, we group statements such that at most one shared variable is accessed and label these ‘atomic’ statements with a number. The labels are chosen identical to the labels in the PVS code, and therefore not completely consecutive.

In every execution step, one of the processes proceeds from one label to a next one. The steps are thus treated as atomic. The atomicity of steps that refer to shared variables more than once is emphasized by enclosing them in angular brackets. Since procedure calls only modify private control data, procedure headers are not always numbered themselves, but their bodies usually have numbered atomic statements.

```

proc find( $a : \text{Address} \setminus \{0\}$ ) : Value =
  local  $r : E\text{Value} ; n, l : \text{Nat} ; h : \text{pointer to Hashtable} ;$ 
5:    $h := H[\text{index}] ; n := 0 ; \{cnt := 0\} ;$ 
6:    $l := h.\text{size} ;$ 
  repeat
7:      $\langle r := h.\text{table}[\text{key}(a, l, n)] ;$ 
       $\{ \text{if } r = \text{null} \vee a = \text{ADR}(r) \text{ then } cnt++ ; (\text{fS}) \text{ end } \} \rangle ;$ 
8:     if  $r = \text{done}$  then
      refresh() ;
10:     $h := H[\text{index}] ; n := 0 ;$ 
11:     $l := h.\text{size} ;$ 
      else  $n++$  end ;
13:  until  $r = \text{null} \vee a = \text{ADR}(r) ;$ 
14:  return val( $r$ ) .

```

In order to prove correctness, we add between braces instructions that only modify auxiliary variables, like the specification variables X and rS and other auxiliary variables to be introduced later. The part between braces is comment for the implementation, it only serves in the proof of correctness. The private auxiliary variable cnt of type Nat counts the number of times (fS) is executed and serves to prove that (fS) is executed precisely once in every call of *find*.

This procedure matches the code of an ordinary find in a hash table with open addressing, except for the code at the condition $r = \text{done}$. This code is needed for the case that the value at address a has been copied, in which case the new table must be located. Locating the new table is carried out by the procedure *refresh*, which is discussed in Section 2.3.5.

In line 7, the accessed hash table should be valid (see invariants *fi4* and *He4* in Appendix A.1). After *refresh* the local variables *n*, *h* and *l* must be reset, to restart the search in the new hash table. If the procedure terminates, the specifying atomic command (fS) has been executed precisely once (see invariant *Cn1*) and the return values of the specification and the implementation are equal (see invariant *Co1*). If the operation succeeds, the return value must be a valid entry currently associated with the given address in the current hash table. It is not evident but it has been proved that the linear search of the process executing *find* cannot be violated by other processes (see invariants *Cu9*, *Cu10* and *fi8*), i.e. no other process can *delete*, *insert*, or *rewrite* an entry associated with the same address (as what the process is looking for) in the region where the process has already searched.

We require that every valid hash table contains at least one entry **null** or **done**. Therefore, the local variable *n* in the procedure *find* never goes beyond the size of the hash table (see invariants *Cu1*, *fi4*, *fi5* and axiom *Ax2*). When the **bound** of the new hash table is tuned properly before use (see invariants *Ne7*, *Ne8*), the hash table will not be updated too frequently, and termination of the procedure *find* can be guaranteed.

Procedure *delete*

To some extent, deletion is similar to finding. Since *r* is a local variable to the procedure *delete*, we regard 18a and 18b as two parts of atomic instruction 18. If the entry is found in the table, then at line 18b this entry is overwritten with the designated element **del**.

```

proc delete(a : Address \ {0}) : Bool =
  local r : EValue ; k, l, n : Nat ;
    h : pointer to Hashtable ; suc : Bool ;
15:   h := H[index] ; suc := false ; {cnt := 0} ;
16:   l := h.size ; n := 0 ;
  repeat
17:     k := key(a, l, n) ;
      ⟨ r := h.table[k] ;
        { if r = null then cnt++ ; (dS) end } ⟩ ;
18a:   if oldp(r) then
      refresh() ;
20:   h := H[index] ;

```

```

21:         l := h.size ; n := 0 ;
        elseif a = ADR(r) then
18b:         ⟨ if r = h.table[k] then
                suc := true ; h.table[k] := del ;
                { cnt++ ; (dS) ; Y[k] := del }
                end ⟩ ;
        else n++ end ;
        until suc ∨ r = null ;
25:         if suc then h.dels++ end ;
26:         return suc .

```

The repetition in this procedure has two ways to terminate. Either deletion fails with $r = \mathbf{null}$ in 17, or deletion succeeds with $r = h.\mathbf{table}[k]$ in 18b. In the latter case, we have in one atomic statement a double access of the shared variable $h.\mathbf{table}[k]$. This is a so-called compare&swap instruction. Atomicity is needed here to preclude interference. The specifying command (dS) is executed either in 17 or in 18b, and it is executed precisely once (see invariant *Cn2*), since in 18b the guard $a = \mathit{ADR}(r)$ implies $r \neq \mathbf{null}$ (see invariant *de1* and axiom *Ax1*).

In order to remember the address from the value rewritten to **done** after the value is being copied in the procedure *moveContents*, in 18, we introduce a new auxiliary shared variable Y of type array of *EValue*, whose contents equals the corresponding contents of the current hash table almost everywhere except that the values it contains are not tagged as *old* or rewritten as **done** (see invariants *Cu9*, *Cu10*).

Since we postpone the increment of $h.\mathbf{dels}$ until line 25, the field \mathbf{dels} is a lower bound of the number of positions deleted in the hash table (see invariant *Cu4*).

Procedure *insert*

The procedure for insertion in the table is given below. Basically, it is the standard algorithm for insertion in a hash table with open addressing. Notable is line 28 where the current process finds that the current hash table too full, and orders a new table to be made. We assume that $h.\mathbf{bound}$ is a number less than $h.\mathbf{size}$ (see invariant *Cu3*), which is tuned for optimal performance.

Furthermore, in line 35, it can be detected that values in the hash table have been

marked *old*, which is a sign that hash table h is outdated, and the new hash table must be located to perform the insertion.

```

proc insert( $v : Value \setminus \{\mathbf{null}\}$ ) : Bool =
  local  $r : EValue ; k, l, n : Nat ; h : \mathbf{pointer\ to\ Hashtable} ;$ 
     $suc : Bool ; a : Address := ADR(v) ;$ 
27:    $h := H[index] ; \{cnt := 0\} ;$ 
28:   if  $h.occ > h.bound$  then
    newTable() ;
30:    $h := H[index]$  end ;
31:    $n := 0 ; l := h.size ; suc := false ;$ 
  repeat
32:      $k := key(a, l, n) ;$ 
33:      $\langle r := h.table[k] ;$ 
       $\{ \mathbf{if} \ a = ADR(r) \ \mathbf{then} \ cnt++ ; (iS) \ \mathbf{end} \} \rangle ;$ 
35a:    if  $oldp(r)$  then
      refresh() ;
36:     $h := H[index] ;$ 
37:     $n := 0 ; l := h.size ;$ 
    elsif  $r = \mathbf{null}$  then
35b:     $\langle \mathbf{if} \ h.table[k] = \mathbf{null} \ \mathbf{then}$ 
       $suc := true ; h.table[k] := v ;$ 
       $\{ cnt++ ; (iS) ; Y[k] := v \}$ 
      end  $\rangle ;$ 
    else  $n++$  end ;
  until  $suc \vee a = ADR(r) ;$ 
41:   if  $suc$  then  $h.occ++$  end ;
42:   return  $suc$  .

```

Instruction 35b is a version of compare&swap. Procedure *insert* terminates successfully when the insertion to an empty slot is completed, or it fails when there already exists an entry with the given address currently in the hash table (see invariant *Co3* and the specification of *insert*).

Procedure assign

Procedure *assign* is almost the same as *insert* except that it rewrites an entry with a given value even when the associated address is not empty. We provide it without further comments.

```

proc assign( $v : Value \setminus \{\mathbf{null}\}$ ) =
  local  $r : EValue ; k, l, n : Nat ; h : \mathbf{pointer\ to\ Hashtable} ;$ 
     $suc : Bool ; a : Address := ADR(v) ;$ 
43:  $h := H[index] ; \{cnt := 0\} ;$ 
44: if  $h.occ > h.bound$  then
    newTable() ;
46:  $h := H[index]$  end ;
47:  $n := 0 ; l := h.size ; suc := false ;$ 
  repeat
48:  $k := key(a, l, n) ;$ 
49:  $r := h.table[k] ;$ 
50a: if  $oldp(r)$  then
    refresh() ;
51:  $h := H[index] ;$ 
52:  $n := 0 ; l := h.size ;$ 
    elseif  $r = \mathbf{null} \vee a = ADR(r)$  then
50b:  $\langle$  if  $h.table[k] = r$  then
     $suc := true ; h.table[k] := v ;$ 
     $\{ cnt++ ; (aS) ; Y[k] := v \}$ 
    end  $\rangle$ 
    else  $n++$  end ;
  until  $suc$  ;
57: if  $r = \mathbf{null}$  then  $h.occ++$  end ;
end.

```

2.3.5 Memory management and concurrent migration

In this section, we provide the public procedures *getAccess* and *releaseAccess* and the auxiliary procedures *refresh* and *newTable* which are responsible for allocation and deallocation.

We begin with the treatment of memory by providing a model of the heap.

The model of the heap

We *model* the **Heap** as an infinite array of hash tables, declared and initialized in the following way:

```

Heap : array Nat of Hashtable := ([Nat] $\perp$ ) ;
H_index : Nat := 1 .

```

So, initially, $\text{Heap}[i] = \perp$ for all indices i . The indices of array **Heap** are the pointers to hash tables. We thus simply regard **pointer to Hashtable** as a synonym of *Nat*. Therefore, the notation $h.\text{table}$ used elsewhere in this chapter stands for $\text{Heap}[h].\text{table}$. Since we reserve 0 (to be distinguished from the absent hash table \perp and the absent value **null**) for the null pointer (i.e. $\text{Heap}[0] = \perp$, see invariant *HeI*), we initialize **H_index**, which is the index of the next hash table, to be 2 instead of 0 or 1. Allocation of memory is modeled in

```

proc allocate(s, b : Nat) : Nat =
  ⟨ Heap[H_index] := blank hash table with size = s, bound = b,
    occ = dels = 0 ;
    H_index++ ⟩ ;
return H_index ;

```

We assume that *allocate* sets all values in the hash table $\text{Heap}[\text{H_index}]$ to **null**, and also sets its fields **size** and **bound** as specified. The variables **occ** and **dels** are set to 0 because the hash table is completely filled with the value **null**.

Deallocation of hash tables is modeled by

```

proc deAlloc(h : Nat) =
  ⟨ assert  $\text{Heap}[h] \neq \perp$  ;  $\text{Heap}[h] := \perp$  ⟩
end .

```

The **assert** here indicates the obligation to prove that *deAlloc* is called only for allocated memory.

Procedure *getAccess*

The procedure *getAccess* is defined as follows.

```

proc getAccess() =
  loop
59:   index := currInd;
60:   prot[index]++ ;
61:   if index = currInd then
62:     busy[index]++ ;
63:     if index = currInd then return ;
        else releaseAccess(index) end ;
65:   else prot[index]-- end ;
  end
end.

```

This procedure is a bit tricky. When the process reaches line 62, the *index* has been protected not to be used for creating a new hash table in the procedure *newTable* (see invariants *pr2*, *pr3* and *nT12*).

The hash table pointer $H[index]$ must contain the valid contents after the procedure *getAccess* returns (see invariants *Ot3*, *He4*). So, in line 62, **busy** is increased, guaranteeing that the hash table will not inadvertently be destroyed (see invariant *bu1* and line 69). Line 63 needs to check the *index* again in case that instruction 62 has the precondition that the hash table is not valid. Once some process gets hold of one hash table after calling *getAccess*, no process can throw it away until the process releases it (see invariant *rA7*).

Procedure *releaseAccess*

The procedure *releaseAccess* is given by

```

proc releaseAccess(i : 1 .. 2P) =
  local h : pointer to Hashtable ;
67:   h := H[i] ;
68:   busy[i]-- ;
69:   if h ≠ 0 ∧ busy[i] = 0 then
70:     ⟨ if H[i] = h then H[i] := 0 ; ⟩
71:     deAlloc(h) ;
        end ;
  end ;

```

```

72:     prot[i]--;
      end.

```

The test $h \neq 0$ at 69 is necessary since it is possible that $h = 0$ at the lines 68 and 69. This occurs e.g. in the following scenario. Assume that process p is at line 62 with $index \neq currInd$, while the number $i = index$ satisfies $H[i] = 0$ and $busy[i] = 0$. Then process p increments $busy[i]$, calls $releaseAccess(i)$, and arrives at 68 with $h = 0$.

Since $deAlloc$ in line 71 accesses a shared variable, we have separated its call from 70. The counter $busy[i]$ is used to protect the hash table from premature deallocation. Only if $busy[i]=0$, $H[i]$ can be released. The main problem of the design at this point is that it can happen that several processes concurrently execute $releaseAccess$ for the same value of i , with interleaving just after the decrement of $busy[i]$. Then they all may find $busy[i] = 0$. Therefore, a bigger atomic command is needed to ensure that precisely one of them sets $H[i]$ to 0 (line 70) and calls $deAlloc$. Indeed, in line 71, $deAlloc$ is called only for allocated memory (see invariant $rA3$). The counter $prot[i]$ can be decreased since position i is no longer used by this process.

Procedure *newTable*

When the current hash table has been used for some time, some actions of the processes may require replacement of this hash table. Procedure *newTable* is called when the number of occupied positions in the current hash table exceeds the bound (see lines 28, 44). Procedure *newTable* tries to allocate a new hash table as the successor of the current one. If several processes call *newTable* concurrently, they need to reach consensus on the choice of an index for the next hash table (in line 84). A newly allocated hash table that will not be used must be deallocated again.

```

      proc newTable() =
        local i : 1 .. 2P ; b, bb : Bool ;
77:     while next[index] = 0 do
78:         choose i ∈ 1 .. 2P ;
           ⟨ b := (prot[i] = 0) ;
             if b then prot[i] := 1 end ⟩ ;
           if b then
81:             busy[i] := 1 ;

```

```

82:      choose bound > H[index].bound - H[index].dels + 2P ;
      choose size > bound + 2P ;
      H[i] := allocate(size, bound) ;
83:      next[i] := 0 ;
84:      ⟨ bb := (next[index] = 0) ;
      if bb then next[index] := i end ⟩ ;
      if -bb then releaseAccess(i) end ;
      end end ;
      refresh() ;
end .

```

In command 82, we allocate a new blank hash table (see invariant *nT8*), of which the **bound** is set greater than $H[index].\text{bound} - H[index].\text{dels} + 2P$ in order to avoid creating a too small hash table (see invariants *nT6*, *nT7*).

We require the **size** of a hash table to be more than $\text{bound} + 2P$ because of the following scenario: P processes find “ $h.\text{occ} > h.\text{bound}$ ” at line 28 and call *newtable*, *refresh*, *migrate*, *moveContents* and *moveElement* one after the other. After moving some elements, all processes but process p sleep at line 126 with $b_{mE} = \text{true}$ (b_{mE} is the local variable b of procedure *moveElement*). Process p continues the migration and updates the new current index when the migration completes. Then, process p does several insertions to let the **occ** of the current hash table reach one more than its **bound**. Just at that moment, $P - 1$ processes wake up, increase the **occ** of the current hash table to be $P - 1$ more, and return to line 30. Since $P - 1$ processes insert different values in the hash table, after $P - 1$ processes finish their insertions, the **occ** of the current hash table reaches $2P - 1$ more than its **bound**.

It may be useful to make **size** larger than $\text{bound} + 2P$ to avoid too many collisions, e.g. with a constraint $\text{size} \geq \alpha \cdot \text{bound}$ for some $\alpha > 1$. If we did not introduce **dels**, every migration would force the sizes to grow, so that our hash table would require unbounded space for unbounded life time. We introduced **dels** to avoid this.

Strictly speaking, instruction 82 inspects one shared variable, $H[index]$, and modifies three other shared variables, viz. $H[i]$, $\text{Heap}[H_index]$, and H_index . In general, we split such multiple shared variable accesses in separate atomic commands. Here the accumulation is harmless, since the only possible interferences are with other allocations at line 82 and deallocations at line 71. In view of the invariant *Ha2*, all deallocations are at pointers $h < H_index$. Allocations do not interfere because they contain the increment $H_index++$

(see procedure *allocate*).

The procedure *newTable* first searches for a free index i , say by round robin. We use a nondeterministic choice. Once a free index has been found, a hash table is allocated and the index gets an indirection to the allocated address. Then the current index gets a **next** pointer to the new index, unless this pointer has been set already.

The variables `prot[i]` are used primarily as counters with atomic increments and decrements. In 78, however, we use an atomic test-and-set instruction. Indeed, separation of this instruction in two atomic instructions is incorrect, since that would allow two processes to grab the same index i concurrently.

Procedure *migrate*

After the choice of the new hash table, the procedure *migrate* serves to transfer the contents in the current hash table to the new hash table by calling a procedure *moveContents* and to update the current hash table pointer afterwards. Migration is complete when at least one of the (parallel) calls to *migrate* has terminated.

```

proc migrate() =
    local  $i : 0 \dots 2P$ ;  $h$  : pointer to Hashtable ;  $b$  : Bool ;
94:    $i := \text{next}[\text{index}]$ ;
95:    $\text{prot}[i]++$  ;
97:   if  $\text{index} \neq \text{currInd}$  then
98:      $\text{prot}[i]--$  ;
    else
99:      $\text{busy}[i]++$  ;
100:     $h := \text{H}[i]$  ;
101:    if  $\text{index} = \text{currInd}$  then
         $\text{moveContents}(\text{H}[\text{index}], h)$  ;
103:     $\langle b := (\text{currInd} = \text{index})$  ;
        if  $b$  then  $\text{currInd} := i$  ;  $\{Y := \text{H}[i].\text{table}\}$ 
        end  $\rangle$  ;
    if  $b$  then
104:         $\text{busy}[\text{index}]--$  ;
105:         $\text{prot}[\text{index}]--$  ;

```

```

    end end ;
    releaseAccess(i) ;
end end .

```

According to invariants *mi4* and *mi5*, it is an invariant that $i = \text{next}(\text{index}) \neq 0$ holds after instruction 94.

Line 103 contains a compare&swap instruction to update the current hash table pointer when some process finds that the migration is finished while `currInd` is still identical to its *index*, which means that *i* is still used for the next current hash table (see invariant *mi5*). The increments of `prot[i]` and `busy[i]` here are needed to protect the next hash table. The decrements serve to avoid memory loss.

Procedure *refresh*

In order to avoid that a delayed process starts migration of an old hash table, we encapsulate *migrate* in *refresh* in the following way.

```

proc refresh() =
90:   if index ≠ currInd then
        releaseAccess(index) ;
        getAccess() ;
      else migrate() end ;
end.

```

When *index* is outdated, the process needs to call *releaseAccess* to abandon its hash table and *getAccess* to acquire the present pointer to the current hash table. Otherwise, the process can join the migration.

Procedure *moveContents*

Procedure *moveContents* has to move the contents of the current table to the next current table. All processes that have access to the table, may also participate in this migration. Indeed, they cannot yet use the new table (see invariants *Ne1* and *Ne3*). We have to take care that delayed actions on the current table and the new table are carried out or abandoned correctly (see invariants *Cu1* and *mE10*). Migration requires that every value in the current table be moved to a unique position in the new table (see invariant *Ne19*).

Procedure *moveContents* uses a private variable *toBeMoved* that ranges over sets of locations. The procedure is given by

```

proc moveContents(from, to : pointer to Hashtable) =
  local i : Nat ; b : Bool ; v : EValue ; toBeMoved : set of Nat ;
  toBeMoved := {0, ..., from.size - 1} ;
110: while currInd = index ∧ toBeMoved ≠ ∅ do
111:   choose i ∈ toBeMoved ;
      v := from.table[i] ;
      if v = done then
112:       toBeMoved := toBeMoved - {i} ;
      else
114:       ⟨ b := (v = from.table[i]) ;
          if b then from.table[i] := old(val(v)) end ⟩ ;
          if b then
116:           if val(v) ≠ null then moveElement(val(v), to) end ;
117:           from.table[i] := done ;
118:           toBeMoved := toBeMoved - {i} ;
          end end end ;
  end .

```

Note that the value is tagged as outdated before it is copied (see invariant *mC11*). After tagging, the value cannot be deleted or assigned until the migration has been completed. Tagging must be done atomically, since otherwise an interleaving deletion may be lost. When indeed the value has been copied to the new hash table, it becomes **done** in the old hash table in line 117. This has the effect that other processes need not wait for this process to complete procedure *moveElement*, but can help with the migration of this value if needed.

Since the address is lost after being rewritten to **done**, we had to introduce the shared auxiliary hash table *Y* to remember its value for the proof of correctness. This could have been avoided by introducing a second tagging bit, say for “very old”.

The processes involved in the same migration should not use the same strategy for choosing *i* in line 111, since it is advantageous that *moveElement* is called often with different values. They may exchange information: any of them may replace its set *toBeMoved* by the

intersection of that set with the set *toBeMoved* of another one. We do not give a preferred strategy here, one can refer to algorithms for the *write-all* problem [24, 46].

Procedure *moveElement*

The procedure *moveElement* moves a value to the new hash table. Note that the value is tagged as outdated in *moveContents* before *moveElement* is called.

```

proc moveElement(v : Value \ {null}, to : pointer to Hashtable) =
  local a : Address ; k, m, n : Nat ; w : EValue ; b : Bool ;
120:  n := 0 ; b := false ; a := ADR(v) ; m := to.size ;
  repeat
121:    k := key(a, m, n) ; w := to.table[k] ;
    if w = null then
123:      ⟨ b := (to.table[k] = null) ;
        if b then to.table[k] := v end ⟩ ;
      else n++ end ;
125:  until b ∨ a = ADR(w) ∨ currInd ≠ index ;
126:  if b then to.occ++ end
  end .

```

The value is only allowed to be inserted once in the new hash table (see invariant *Ne19*), since otherwise the main property of open addressing would be violated. In total, four situations can occur in the procedure *moveElement*:

- the current location *k* contains a value with a different address. The process increases *n* to inspect the next location.
- the current location *k* contains a value with the same address. This means that the value has already been copied to the new hash table, the process therefore terminates.
- the current location *k* is an empty slot. The process inserts *v* and returns. If insertion fails, since another process filled the empty slot in between, the search is continued.
- when *index* happens to differ from *currInd*, the entire migration has been completed.

While the current hash table pointer is not updated yet, there exists at least one **null** entry in the new hash table (see invariants *Ne8*, *Ne22* and *Ne23*), hence the local variable *n*

in the procedure *moveElement* never goes beyond the size of the hash table (see invariants *mE3* and *mE8*), and the termination is thus guaranteed.

2.4 Correctness (Safety)

In this section, we describe the proof of safety of the algorithm. The main aspects of safety are functional correctness, atomicity, and absence of memory loss. These aspects are formalized in eight invariants described in section 2.4.1. To prove these invariants, we need many other invariants. These are listed in Appendix A.1. In section 2.4.2, we sketch the verification of some of the invariants by informal means. In section 2.4.3, we describe how the theorem prover PVS is used in the verification. As exemplified in 2.4.2, Appendix A.2 gives the dependencies between the invariants.

Notational Conventions. Recall that there are at most P processes with process identifiers ranging from 1 up to P . We use p, q, r to range over process identifiers, with a preference for p . Since the same program is executed by all processes, every private variable name of a process $\neq p$ is extended with the suffix “.” + “process identifier”. We do not do this for process p . So, e.g., the value of a private variable x of process q is denoted by $x.q$, but the value of x of process p is just denoted by x . In particular, $pc.q$ is the program location of process q . It ranges over all integer labels used in the implementation.

When local variables in different procedures have the same names, we add an abbreviation of the procedure name as a subscript to the name. We use the following abbreviations: *fi* for *find*, *del* for *delete*, *ins* for *insert*, *ass* for *assign*, *gA* for *getAccess*, *rA* for *releaseAccess*, *nT* for *newTable*, *mig* for *migrate*, *ref* for *refresh*, *mC* for *moveContents*, *mE* for *moveElement*.

In the implementation, there are several places where the same procedure is called, say *getAccess*, *releaseAccess*, etc. We introduce auxiliary private variables *return*, local to such a procedure, to hold the return location. We add a procedure subscript to distinguish these variables according to the above convention.

If V is a set, $\#V$ denotes the number of elements of V . If b is a boolean, then $\#b = 0$ when b is false, and $\#b = 1$ when b is true. Unless explicitly defined otherwise, we always (implicitly) universally quantify over addresses a , values v , non-negative integer numbers k, m , and n , natural number l , processes p, q and r . Indices i and j range over $[1, 2P]$. We abbreviate $H(\text{currInd}).\text{size}$ as *curSize*.

In order to avoid using too many parentheses, we use the usual binding order for the operators. We give “ \wedge ” higher priority than “ \vee ”. We use parentheses whenever necessary.

2.4.1 Main properties

We have proved the following three safety properties of the algorithm. Firstly, the access procedures *find*, *delete*, *insert*, *assign*, are functionally correct. Secondly they are executed atomically. The third safety property is absence of memory loss.

Functional correctness of *find*, *delete*, *insert* is the condition that the result of the implementation is the same as the result of the specification (fS), (dS), (iS). This is expressed by the required invariants:

$$\text{Co1: } pc = 14 \Rightarrow \text{val}(r_{fi}) = rS_{fi}$$

$$\text{Co2: } pc \in \{25, 26\} \Rightarrow \text{suc}_{del} = \text{suc}S_{del}$$

$$\text{Co3: } pc \in \{41, 42\} \Rightarrow \text{suc}_{ins} = \text{suc}S_{ins}$$

Note that functional correctness of *assign* holds trivially since it does not return a result.

According to the definition of atomicity in chapter 13 of [52], atomicity means that each execution of one of the access procedures contains precisely one execution of the corresponding specifying action (fS), (dS), (iS), (aS). We introduced the private auxiliary variables *cnt* to count the number of times the specifying action is executed. Therefore, atomicity is expressed by the invariants:

$$\text{Cn1: } pc = 14 \Rightarrow \text{cnt}_{fi} = 1$$

$$\text{Cn2: } pc \in \{25, 26\} \Rightarrow \text{cnt}_{del} = 1$$

$$\text{Cn3: } pc \in \{41, 42\} \Rightarrow \text{cnt}_{ins} = 1$$

$$\text{Cn4: } pc = 57 \Rightarrow \text{cnt}_{ass} = 1$$

We interpret absence of memory loss to mean that the number of allocated hash tables is bounded. More precisely, we prove that this number is bounded by $2P$. This is formalized in the invariant:

$$\text{No1: } \#\{k \mid k < \text{H_index} \wedge \text{Heap}(k) \neq \perp\} \leq 2P$$

An important safety property is that no process accesses deallocated memory. Since most procedures perform memory accesses, by means of pointers that are local variables, the proof of this is based on a number of different invariants. Although this is not explicit

in the specification, it has been checked because the theorem prover PVS does not allow access to deallocated memory as this would violate type correctness conditions.

2.4.2 Intuitive proof

The eight correctness properties (invariants) mentioned above have been completely proved with the interactive proof checker of PVS. The use of PVS did not only take care of the delicate bookkeeping involved in the proof, it could also deal with many trivial cases automatically. At several occasions where PVS refused to let a proof be finished, we actually found a mistake and had to correct previous versions of this algorithm.

In order to give some feeling for the proof, we describe some proofs. For the complete mechanical proof, we refer the reader to [33]. Note that, for simplicity, we assume that all non-specific private variables in the proposed assertions belong to the general process p , and general process q is an active process that tries to threaten some assertion (p may equal q).

Proof of invariant *Co1* (as claimed in 2.4.1). According to Appendix A.2, the stability of *Co1* follows from the invariants *Ot3*, *fi1*, *fi10*, which are given in Appendix A.1. Indeed, *Ot3* implies that no procedure returns to location 14. Therefore all return statements falsify the antecedent of *Co1* and thus preserve *Co1*. Since r_{f_i} and rS_{f_i} are private variables to process p , *Co1* can only be violated by process p itself (establishing *pc at 14*) when p executes 13 with $r_{f_i} = \mathbf{null} \vee a_{f_i} = \mathit{ADR}(r_{f_i})$. This condition is abbreviated as $\mathit{Find}(r_{f_i}, a_{f_i})$. Invariant *fi10* then implies that action 13 has the precondition $\mathit{val}(r_{f_i}) = rS_{f_i}$, so then it does not violate *Co1*. In PVS, we used a slightly different definition of Find , and we applied invariant *fi1* to exclude that r_{f_i} is **done** or **del**, though invariant *fi1* is superfluous in this intuitive proof. \square

Proof of invariant *Ot3*. Since the procedures $\mathit{getAccess}$, $\mathit{releaseAccess}$, $\mathit{refresh}$, $\mathit{newTable}$ are called only at specific locations in the algorithm, it is easy to list the potential return addresses. Since the variables *return* are private to process p , they are not modified by other processes. Stability of *Ot3* follows from this. As we saw in the previous proof, *Ot3* is used to guarantee that no unexpected jumps occur. \square

Proof of invariant *fi10*. According to Appendix A.2, we only need to use *fi9* and *Ot3*. Let us use the abbreviation $k = \mathit{key}(a_{f_i}, l_{f_i}, n_{f_i})$. Since r_{f_i} and rS_{f_i} are both private variables,

they can only be modified by process p when p is executing statement 7. We split this situation into two cases

1. with precondition $Find(h_{f_i}.table[k], a_{f_i})$
 After execution of statement 7, r_{f_i} becomes $h_{f_i}.table[k]$, and rS_{f_i} becomes $X(a_{f_i})$. By $fi9$, we get $val(r_{f_i}) = rS_{f_i}$. Therefore the validity of $fi10$ is preserved.
2. otherwise.
 After execution of statement 7, r_{f_i} becomes $h_{f_i}.table[k]$, which then falsifies the antecedent of $fi10$. □

Proof of invariant $fi9$. According to Appendix A.2, we proved that $fi9$ follows from $Ax2$, $fi1$, $fi3$, $fi4$, $fi5$, $fi8$, $Ha4$, $He4$, $Cu1$, $Cu9$, $Cu10$, and $Cu11$. We abbreviate $key(a_{f_i}, l_{f_i}, n_{f_i})$ as k . We deduce $h_{f_i} = H(index)$ from $fi4$, $H(index)$ is not \perp from $He4$, and k is below $H(index).size$ from $Ax2$, $fi4$ and $fi3$. We split the proof into two cases:

1. $index \neq currInd$: By $Ha4$, it follows that $H(index) \neq H(currInd)$. Hence from $Cu1$, we obtain $h_{f_i}.table[k] = done$, which falsifies the antecedent of $fi9$.
2. $index = currInd$: By premise $Find(h_{f_i}.table[k], a_{f_i})$, we know that $h_{f_i}.table[k] \neq done$ because of $fi1$. By $Cu9$ and $Cu10$, we obtain $val(h_{f_i}.table[k]) = val(Y[k])$. Hence it follows that $Find(Y[k], a_{f_i})$. Using $fi8$, we obtain

$$\forall m < n_{f_i} : \neg Find(Y[key(a_{f_i}, curSize, m)], a_{f_i})$$

We get n_{f_i} is below $curSize$ because of $fi5$. By $Cu11$, we conclude

$$X(a_{f_i}) = val(h_{f_i}.table[k])$$

□

2.4.3 The model in PVS

Our proof architecture (for one property) can be described as a dynamically growing tree in which each node is associated with an assertion. We start from a tree containing only one node, the proof goal, which characterizes some property of the system. We expand the tree

by adding some new children via proper analysis of an unproved node (top-down approach, which requires a good understanding of the system). The validity of that unproved node is then reduced to the validity of its children and the validity of some less or equally deep nodes.

Normally, simple properties of the system are proved with appropriate precedence, and then used to help establish more complex ones. It is not a bad thing that some property that was taken for granted turns out to be not valid. Indeed, it may uncover a defect of the algorithm, but in any case it leads to new insights in it.

We model the algorithm as a transition system [53], which is described in the language of PVS in the following way. As usual in PVS, states are represented by a record with a number of fields:

```

State : TYPE = [#
% global variables
...
  busy : [ range(2*P) → nat ],
  prot : [ range(2*P) → nat ],
...
% private variables:
  index : [ range(P) → range(2*P) ],
...
  pc : [ range(P) → nat ], % private program counters
...
% local variables of procedures, also private to each process:
% find
  a_find : [ range(P) → Address ],
  r_find : [ range(P) → EValue ],
...
% getAccess
  return_getAccess : [ range(P) → nat ],
...
#]
```

where $range(P)$ stands for the range of integers from 1 to P .

Note that private variables are given with as argument a process identifier. Local variables are distinguished by adding their procedure's names as suffixes.

An action is a binary relation on states: it relates the state prior to the action to the

state following the action. The system performed by a particular process is then specified by defining the precondition of each action as a predicate on the state and also the effect of each action in terms of a state transition. For example, line 5 of the algorithm is described in PVS as follows:

```
% corresponding to statement find5: h := H[index]; n := 0;
find5(i,s1,s2) : bool =
  pc(s1)(i)=5 AND
  s2 = s1 WITH [ (pc)(i) := 6,
                 (n_find)(i) := 0,
                 (h_find)(i) := H(s1)(index(s1)(i)) ]
```

where i is a process identifier, $s1$ is a pre-state, $s2$ is a post-state.

Since our algorithm is concurrent, the global transition relation is defined as the disjunction of all atomic actions.

```
% transition steps
step(i,s1,s2) : bool =
  find5(i,s1,s2) or find6(i,s1,s2) or ...
  delete15(i,s1,s2) or delete16(i,s1,s2) or ...
  ...
```

Stability for each invariant is proved by a PVS *Theorem* of the form:

```
% Theorem about the stability of invariant fi10
IV_fi10: THEOREM
  forall (u,v : state, q : range(P) ) :
    step(q,u,v) AND fi10(u) AND fi9(u) AND ot3(u)
    => fi10(v)
```

To ensure that all proposed invariants are stable, there is a global invariant *INV*, which is the conjunction of all proposed invariants.

```
% global invariant
INV(s:state) : bool =
  He3(s) and He4(s) and Cu1(s) and ...
  ...
```

```
% Theorem about the stability of the global invariant INV
IV_INV: THEOREM
  forall (u,v : state, q : range(P) ) :
    step(q,u,v) AND INV(u) => INV(v)
```

We define `Init` as all possible initial states, for which all invariants must be valid.

```
% initial state
Init: { s : state |
  (forall (p: range(P)):
    pc(s)(p)=0 and ...
    ...) and
  (forall (a: Address):
    X(s)(a)=null) and
  ...
}

% The initial condition can be satisfied by the global invariant INV
IV_Init: THEOREM
  INV(Init)
```

The PVS code contains eleven preconditions to imply well-definedness: e.g. in *find7*, the hash table must be non-NIL and ℓ must be its size.

```
% corresponding to statement find7
find7(i,s1,s2) : bool =
  i?(Heap(s1)(h_find(s1)(i))) and
  l_find(s1)(i)=size(i_(Heap(s1)(h_find(s1)(i)))) and
  pc(s1)(i)=7 and
  ...
```

All preconditions are allowed, since we can prove lock-freedom in the following form. In every state $s1$ that satisfies the global invariant, every process q can perform a step, i.e., there is a state $s2$ with $(s1, s2) \in \text{step}$ and $pc(s1, q) \neq pc(s2, q)$. This is expressed in PVS by

```
% theorem for lock-freedom
IV_prog: THEOREM
  forall (u: state, q: range(P) ) :
    INV(u) => exists (v: state): pc(u)(q) /= pc(v)(q) and step(q,u,v)
```

2.5 Correctness (Progress)

In this section, we prove that our algorithm is lock-free, and that it is wait-free for several subtasks. However, the proof was not checked with PVS.

Recall that an algorithm is called *lock-free* if always at least some process will finish its task in a finite number of steps, regardless of delays or failures by other processes. This means that no process can block the applications of further operations to the data structure, although any particular operation need not terminate since a slow process can be passed infinitely often by faster processes. We say that an operation is *wait-free* if any process involved in that operation is guaranteed to complete it in a finite number of its own steps, regardless of the (in)activity of other processes.

2.5.1 The easy part of progress

It is clear that *releaseAccess* is wait-free. It follows that the wait-freedom of *migrate* depends on wait-freedom of *moveContents*. The loop of *moveContents* is clearly bounded. So, wait-freedom of *moveContents* depends on wait-freedom of *moveElement*. It has been proved that n is bounded by m in *moveElement* (see invariants *mE3* and *mE8*). Since, moreover, $\text{to.table}[k] \neq \mathbf{null}$ is stable, the loop of *moveElement* is also bounded. This concludes the sketch that *migrate* is wait-free.

2.5.2 Progress of newTable

The main part of procedure *newTable* is wait-free. This can be shown informally, as follows. Since we can prove the condition $\text{next}(\text{index}) \neq 0$ is stable while process p stays in the region $[77, 84]$, once the condition $\text{next}(\text{index}) \neq 0$ holds, process p will exit *newTable* in a few rounds.

Otherwise, we may assume that p has precondition $\text{next}(\text{index}) = 0$ before executing line 78. By the invariant

$$\begin{aligned} \text{Ne5: } & pc \in [1, 58] \quad \vee \quad pc \geq 62 \wedge pc \neq 65 \wedge \text{next}(\text{index}) = 0 \\ & \Rightarrow \quad \text{index} = \text{currInd} \end{aligned}$$

we get that $\text{index} = \text{currInd}$ holds and $\text{next}(\text{currInd}) = 0$ from the precondition. We define two sets of integers:

$$\begin{aligned} \text{prSet1}(i) &= \{r \mid \text{index}.r = i \wedge pc.r \notin \{0, 59, 60\}\} \\ \text{prSet2}(i) &= \{r \mid \text{index}.r = i \wedge pc.r \in \{104, 105\} \\ &\quad \vee i_{rA}.r = i \wedge \text{index}.r \neq i \wedge pc.r \in [67, 72] \\ &\quad \vee i_{nT}.r = i \wedge pc.r \in [81, 84] \\ &\quad \vee i_{mig}.r = i \wedge pc.r \geq 97 \} \end{aligned}$$

and consider the sum $\sum_{i=1}^{2P} (\#(prSet1(i)) + \#(prSet2(i)))$. While process p is at line 78, the sum cannot exceed $2P - 1$ because there are only P processes around and process p contributes only once to the sum. It then follows from the pigeon hole principle that there exists $j \in [1, 2P]$ such that $\#(prSet1(j)) + \#(prSet2(j)) = 0$ and $j \neq index.p$. By the invariant

$$pr1: \quad \text{prot}[j] = \#(prSet1(j)) + \#(prSet2(j)) + \#(\text{currInd} = j) \\ + \#(\text{next}(\text{currInd}) = j)$$

we can get that $\text{prot}[j] = 0$ because of $j \neq index.p = \text{currInd}$.

While currInd is constant, no process can modify $\text{prot}[j]$ for $j \neq \text{currInd}$ infinitely often. Therefore, if process p acts infinitely often and chooses its value i in 78 by round robin, process p exits the loop of *newTable* eventually. This shows that the main part of *newTable* is wait-free.

2.5.3 The failure of wait-freedom

Procedure *getAccess* is not wait-free. When the active clients keep changing the current index faster than the new client can observe it, the accessing client is doomed to starvation. In that case, however, the other processes repeatedly succeed. It follows that *getAccess*, *refresh*, and *newTable* are lock-free.

It may be possible to make a queue for the accessing clients which is emptied by a process in *newTable*. The accessing clients must however also be able to enter autonomously. This would at least add another layer of complications. We therefore prefer to treat this failure of wait-freedom as a performance issue that can be dealt with in practice by tuning the sizes of the hash tables.

According to the invariants *fi5*, *de8*, *in8* and *as6*, the primary procedures *find*, *delete*, *insert*, *assign* are loops bounded by $n \leq h.\text{size}$, and n is only reset to 0 during migration. If n is not reset to 0, it is incremented or stays constant. Indeed, the atomic **if** statements in 18b, 35b, and 50b have no **else** parts. In *delete* and *assign*, it is therefore possible that n stays constant without termination of the loop. Since *assign* can modify non-**null** elements of the table, it follows that *delete* and *assign* are not wait-free. This unbounded fruitless activity is possible only when *assign* actions of other processes repeatedly succeed. It follows that the primary procedures are lock-free. This concludes the argument that the system is lock-free.

2.6 Conclusions

Lock-free shared data objects are inherently resilient to halting failures and permit maximum parallelism. We have presented a new practical, lock-free algorithm for concurrently accessible hash tables, which promises more robust performance and reliability than a conventional lock-based implementation. Moreover, the new algorithm is dynamic in the sense that it allows the hash table to grow and shrink as needed.

The algorithm scales up linearly with the number of processes, provided the function *key* and the selection of *i* in line 111 are defined well. This is confirmed by some experiments where random values were stored, retrieved and deleted from the hash table. These experiments indicated that 10^6 insertions, deletions and finds per second and per processor are possible on an SGI powerchallenge with 250Mhz R12000 processors. This figure should only be taken as a rough indicator, since the performance of parallel processing is very much influenced by the machine architecture, the relative sizes of data structures compared to sizes of caches, and even the scheduling of processes on processors.

The correctness proof for our algorithm is noteworthy because of the extreme effort it took to finish it. Formal deduction by human-guided theorem proving can, in principle, verify any correct design, but doing so may require huge amounts of effort, time, or skill. Though PVS provided great help for managing and reusing the proofs, we have to admit that the verification for our algorithm was very complicated due to the complexity of our algorithm. The total verification effort can roughly be estimated to consist of two man years excluding the effort in determining the algorithm and writing the documentation. The whole proof contains around 200 invariants. It takes a 1 Ghz Pentium IV computer around two days to re-run an individual proof for one of the biggest invariants. Without suitable tool support like PVS, we even doubt if it would be possible to complete a reliable proof of such size and complexity.

It may well be possible to simplify the proof and reduce the number of invariants slightly, but we did not work on this. The complete version of the PVS specifications and the whole proof scripts can be found at [33]. Note that we simplified some definitions in this chapter for the sake of presentation.