

The Derivation of Distributed Termination Detection Algorithms from Garbage Collection Schemes

Gerard Tel*

*Department of Computer Science, University of Utrecht,
P.O. Box 80.089, 3508 TB Utrecht, The Netherlands*

Email: gerard@cs.ruu.nl

Friedemann Mattern

*Department of Computer Science, University of Saarland,
Im Stadtwald 36, D 6600 Saarbrücken, Fed. Rep. Germany*

Email: mattern@cs.uni-sb.de

September 1990

Revised September 1991

Abstract

It is shown that the termination detection problem for distributed computations can be modeled as an instance of the garbage collection problem. Consequently, algorithms for the termination detection problem are obtained by applying transformations to garbage collection algorithms. The transformation can be applied to collectors of the “mark-and-sweep” type as well as to reference counting garbage collectors. As examples, the scheme is used to transform the distributed reference counting protocol of Lermen and Maurer, the weighted reference counting protocol, the local reference counting protocol, and Ben-Ari’s mark-and-sweep collector into termination detection algorithms. Known termination detection algorithms as well as new variants are obtained.

CATEGORIES AND SUBJECT DESCRIPTORS: D.1.3 [**Programming Techniques**]: Concurrent Programming; C.2.4 [**Computer-Communication Networks**]: Distributed Systems—*Network Operating Systems*; D.4.2 [**Operating Systems**]: Storage Management—*Distributed Memories*.

GENERAL TERMS: Algorithms, Design, Theory, Verification.

ADDITIONAL KEYWORDS AND PHRASES: Distributed Algorithms, Distributed Termination Detection, Garbage Collection, Program Transformations.

*The work of this author is supported by the ESPRIT II Basic Research Actions Program of the EC under contract no. 3075 (project ALCOM).

1 Introduction

A substantial amount of the research efforts in distributed algorithm design has been devoted to the problem of detecting when a distributed computation has terminated. There are several reasons for the impressive number of publications on this subject. First, because the problem has shown up under varying model assumptions and because there are several solutions for each model, a very large number of different algorithms has emerged. All these algorithms were published separately, as unifying approaches treating a number of algorithms as a class have been rare. Second, the problem of termination detection, being sufficiently easy to define and yet non-trivial to solve, has been seen as a good candidate to illustrate the merits of design or proof methods for distributed algorithms. Third, it has been observed that the fundamental difficulties of the termination detection problem are the same as those of other important problems in distributed computing. Indeed, termination detection algorithms are related to algorithms for computing distributed snapshots [3, 19], detecting deadlocks [5, 23], and approximating a distributed infimum [31, 20]. Thus the problem is seen to be important both from a practical, algorithmical, and from a theoretical, methodological point of view.

From both points of view we consider it useful to recognize general design paradigms for distributed termination detection algorithms. One such paradigm was described by Tel [32]. A new paradigm is presented in this paper: it is shown that the semantics of the termination detection problem is fully contained in the semantics of the garbage collection problem. As a result, termination detection algorithms are obtained as suitable instantiations of garbage collection algorithms.

In the remainder of this section we first introduce the termination detection problem and the distributed garbage collection problem. Section 2 then describes how the termination detection problem can be formulated as garbage collecting one hypothetical object and derives the algorithmical transformation. Section 3 applies the transformation to four known garbage collection algorithms and discusses the resulting termination detection algorithms. Section 4 contains some additional remarks and comments.

1.1 The Termination Detection Problem

In a distributed system where processes communicate only via messages, in general no process has a consistent and up to date view of the global state. As a result, it is non-trivial to decide whether or not the global state is one in which a distributed computation has terminated. Some processes may have finished their local computations, while others are still executing. But tasks may migrate from one process to another, new tasks may be generated, or the receipt of a message may result in renewed computational activity. As a consequence, finished processes may later be computing again.

In general it is not possible for a process to decide whether it will later

\mathbf{S}_p : { $state_p = active$ }
 send a message $\langle M \rangle$

\mathbf{R}_p : { A basic message has arrived }
 receive message $\langle M \rangle$; $state_p := active$

\mathbf{I}_p : { $state_p = active$ }
 $state_p := passive$

Algorithm 1: The actions of the basic computation.

generate new tasks. Therefore it is always assumed that for each process a local *condition of stability* is defined. While this local condition holds, the process does not send messages (belonging to the computation) to other processes, no new tasks are generated by the process, and no initiative of the process itself falsifies the condition of stability. Only the receipt of a message can do so. It now follows that if a global state is reached in which the condition of stability is satisfied (simultaneously) in every process and no messages are in transit, the computation is terminated. A control computation must be superimposed to detect this situation.

1.1.1 Description of the Problem

The problem is described formally as follows. A collection \mathcal{P} of *processes* is considered, communicating by message passing. A process is either *passive* (if its condition of stability is satisfied) or *active* (if the condition is not satisfied). *Active* processes may send messages, but *passive* processes don't. An *active* process may spontaneously become *passive*, but a *passive* process may become *active* only on receipt of a message.

A full description of the possible actions of processes is given in Algorithm 1. Action \mathbf{S}_p is the sending of a message, action \mathbf{R}_p the receipt of a message, and action \mathbf{I}_p the transition of a process from *active* to *passive*. Throughout the paper it is assumed that each action is executed atomically. An action whose name is subscripted with p takes place in process (or object) p . An assertion between braces (“{” and “}”) is a guard and means that the action can only be executed when the assertion is true. Comments are placed between “(“ and “)“.

Define the *termination condition* as:

No process is *active* and no messages are in transit.

The termination condition is stable: once true, it remains so, because the actions are all disabled if it holds. The problem of termination detection now

is to superimpose on the described *basic* computation a *control* computation which enables one or more of the processes to detect when the termination condition holds. A process detects this by entering a special state *terminated*. The following two criteria specify the correctness of the control algorithm.

T1 *Safety*. If any process is in state *terminated* then the termination condition holds.

T2 *Liveness*. If the termination condition holds, then eventually a process will be in the *terminated* state.

A *passive* process may take part in this control computation, and receiving control messages does not make a *passive* process *active*.

Under varying assumptions about the communication semantics the concise description above still allows different variants of the problem. Originally the problem emerged from a CSP context, where the sending and receipt of a message are synchronized with each other. Thus messages in transit can be ignored; the *termination condition for synchronous communication* simply reads “all processes are passive”. The introduction of asynchronous communication complicates the problem, as somehow it must be verified that the channels are empty. This can be done using special *marker* messages (in a FIFO environment) as by Misra [22], acknowledgements as by Dijkstra and Scholten [11], counting of sent and received messages as by Mattern [17], or by assuming an upper bound on message delay as by Tel [32].

1.1.2 Solutions to the Problem

Several classes of solutions to the termination detection problem are known. The most important ones are those based on *probes* and those based on *acknowledgements*.

Probe-based Algorithms. A *probe* is a distributed activity that “visits” all processes in the network. (It can be implemented by a token circulating on a ring, by an echo mechanism, or in many other ways, see Tel [33].) To detect termination using probes, an attempt is made to maintain a state in which all visited processes are *passive*, and no message is underway to a visited process. A violation of this aim occurs when a non-visited process sends a message to a visited one. Usually, if this happens the current probe is marked as unsuccessful, and after its completion a new probe is initiated. (This marking can be done, for example, by assuming a different “color” for processes that caused the violation and probe messages that report about it.) Termination is detected when a probe completes successfully. The best known example in this class is the algorithm by Dijkstra *et al.* [9], a general treatment is given by Tel [32].

Acknowledgement-based Solutions. In these algorithms all messages of the basic computation are acknowledged, but only after all computational activity resulting from it has ceased. That is, if an *active* process receives a message,

it acknowledges it immediately. If a *passive* process receives a message and becomes *active*, it defers the acknowledgement until it is *passive* again, and has received acknowledgements for all messages it sent during the period of activity. When the initiators of the computation are *passive* and have received an acknowledgement for all basic messages, termination is detected. The best known example in this class is the algorithm of Dijkstra and Scholten [11]; generalizations of this algorithm were given by Shavit and Francez [28] and Chandrasekaran and Venkatesan [8].

1.2 The Distributed Garbage Collection Problem

As our approach for deriving termination detection algorithms is based on solutions to the garbage collection problem, we shall now describe this problem. From a practical point of view, algorithms for the garbage collection problem are important for the storage management of programming languages with dynamic objects. They are also used in the implementation of functional programming languages as these languages operate on directed graphs, represented by memory cells referencing each other through pointers. An account of various garbage collection algorithms for multiprocessors and distributed systems was given by Rudalics [25].

Different models for the problem are found in the literature, here a model based on the communicating objects paradigm is presented which is close to the model of Lermen and Maurer [16]. The advantage of this model is that it abstracts from aspects which are not relevant to our purposes, such as processors, memory cells, and the difference between “local” and “remote” references.

This model differs from the one used by Dijkstra *et al.* [10] in two important respects. First, our model is based on addition and deletion of references, giving rise to a dynamic number of references in each object. The model of [10] assumes that references are only overwritten by other references, giving rise to a constant (*viz.*, 2) number of references per object. Where the latter assumption is justified in the case of extremely simple objects (e.g., Lisp storage cells), the former assumption is justified when more sophisticated objects are considered (as it is the case, for example, in object oriented programming languages with dynamic data structures).

A second difference is the much coarser grain of atomicity that is assumed in our model. As the algorithm of [10] was designed for processes communicating via shared memory, a fine grain of atomicity was desirable to minimize synchronization overhead. In message passing systems (which are assumed in our model) a larger grain of atomicity is permissible.

1.2.1 Description of the Problem

An (object-oriented) distributed system consists of a collection \mathcal{O} of cooperating processes called *objects*. A subset of \mathcal{O} is designated as *root objects*. Objects are able to hold *references* to other objects. These references can be transmitted in

-
- CR_p**: { p is reachable and holds an r -reference }
 send a $\langle \mathbf{cop} \ r \rangle$ message to q
- RC_p**: { A $\langle \mathbf{cop} \ r \rangle$ message has arrived }
 receive the $\langle \mathbf{cop} \ r \rangle$ message ;
 insert the r -reference
- DR_p**: { p holds an r -reference }
 delete the r -reference

Algorithm 2: Reference manipulation by the objects.

messages; see below. A reference to an object r will be called an r -reference. An object r is a *descendant* of q if q holds an r -reference or a message containing an r -reference is in transit to q . An object is *reachable* if it is a root object or a descendant of a reachable object. An object p holding an r -reference may *delete* it, after which p no longer holds this reference. Also, a reachable object p holding an r -reference may *copy* the reference to another object q . It is usually assumed that this happens only when q is also reachable, but this assumption is not necessary for our purposes. To copy an r -reference, p sends the r -reference in a message to q , and q will hold an r -reference after receipt of this message. (Note that only reachable objects may copy references they hold. Since “being reachable” is a non-local and non-stable property, this is a non-trivial requirement imposed on the implementation. Fortunately, this is not a problem for our purpose because—as will become clear later—those objects in our transformation which may copy references are always reachable by definition.) An object can have multiple references to the same target object.

Formally, the allowed actions in this model are as in Algorithm 2. Here **CR_p** is the action by which p initiates the copying of an r -reference to q , **RC_p** is the action by which p inserts an r -reference as a result of copying, and **DR_p** is the action by which p deletes an r -reference. The formal similarity of the actions in Algorithm 1 and Algorithm 2 is remarkable. Obviously, being *active* in the former model corresponds to holding an r -reference in the latter model. This similarity will be investigated further down since it plays a central role in our transformation.

We do not make any assumptions about the message communication system, such as that message communication is synchronous or obeys the FIFO rule. Indeed, such assumptions are not necessary for our transformation—only when the transformation is applied to a garbage collection scheme that relies on such assumptions (see, e.g., Section 3.1), they must be made; then the resulting termination detection algorithm relies on them also.

An object is called *garbage* if it is not reachable. As only references to

reachable objects are copied, a garbage object remains garbage forever. For reasons of memory management it is required that garbage objects are identified and collected. This task is taken care of by a garbage collecting algorithm. The following two criteria define the correctness of a garbage collecting algorithm.

G1 *Safety*. If an object is collected, it is garbage.

G2 *Liveness*. If an object is garbage, it will eventually be collected.

1.2.2 Solutions to the Problem

Many solutions have been proposed to the distributed garbage collection problem, most of which fall into one of two categories: collectors of the *reference counting* type and collectors of the *mark-and-sweep* type. Both types of solutions have been known for over 30 years for classical, non-distributed systems [7, 21].

Reference Counting [2, 14, 15, 16, 24, 36]. Collectors of the first type maintain for each non-root object a count of the number of references in existence to that object. References in other objects as well as references in messages are taken into account. The reference count is incremented when a reference to the object is copied, and decremented when such a reference is deleted. When the count for an object drops to zero, it can be concluded that the object is garbage and consequently the object can be collected.

A group of garbage objects, cyclically referencing each other, cannot be collected by a reference counting algorithm, because no reference count drops to zero. Thus the algorithms do not satisfy the liveness condition G2, but only the weaker condition

G3 If there are no references to a non-root object, it will eventually be collected.

When an object is collected its references are deleted, so that a “chain” of garbage objects will eventually be collected entirely if G3 is satisfied. To obtain an algorithm satisfying G2, usually a supplementary algorithm (typically of the mark-and-sweep type) is used to collect cyclic structures of garbage. In our application, however, cyclic structures of garbage objects do not occur, and a supplementary algorithm is not necessary.

Mark-and-sweep [1, 10, 29, 30]. Collectors of the second type mark all reachable objects as such, starting from the roots and recursively marking all descendants of marked objects. In this way all reachable objects become marked eventually. The design of the marking algorithm is complicated by the possibility that references are inserted and deleted during its operation. The objects in the system must cooperate with the marking algorithm, e.g., by also marking objects when references are installed, copied, or removed. A possible design,

presented by Tel *et al.* [34] consists of an algorithm for the marking proper, upon which a termination detection algorithm is superimposed. When the marking phase is terminated, a sweep through all objects is made in which all unmarked objects are collected. These two phases form one cycle of the collector, and cycles are repeated as long as necessary.

2 Termination Detection Using Garbage Collection

In this section we describe how the distributed termination detection problem in general can be modeled as an instance of the garbage collection problem. As a result, solutions to the termination detection problem can be derived from garbage collection algorithms. First the collection \mathcal{O} of objects used for this purpose is described, as well as the behavior of these objects. Next it is shown that the termination condition is equivalent to one particular object becoming garbage. As a result, termination can be detected by a garbage collection algorithm. Concrete examples of this will be presented in Section 3.

Recall that \mathcal{P} is the set of processes whose termination is to be detected. The collection \mathcal{O} of objects consists of one root object A_p for every process p in \mathcal{P} , and a single *indicator object* Z . Object A_p mimics the behavior of process p as far as the basic computation is concerned (it sends and receives p 's basic messages, and has all the variables p has). Object A_p is called *passive* (*active*) when process p is *passive* (*active*). As A_p is a root object, it is always reachable.

The indicator object Z is not a root object. Its only purpose is to indicate the termination condition with its reachability status by the following equivalence, which will be maintained during execution.

(IND) Z is garbage \Leftrightarrow the termination condition holds.

Theorem 2.1 *IND holds when the following two rules are observed:*

R1 *An object holds a Z -reference if and only if it is active.*

R2 *Each message of the basic computation contains a Z -reference.*

Proof. Z is garbage is equivalent to: Z is not a descendant of any of the A_p . By definition, this means that no A_p holds a Z -reference, and to no A_p a message is in transit containing a Z -reference. By R1 and R2 this is equivalent to: no A_p is *active* and to no A_p a message (of the basic computation) is in transit. This is the definition of the termination condition. \square

It must be shown how R1 and R2 can be maintained. It is possible to ensure through proper initialization that R1 and R2 hold initially. To this end, assume that *active* objects are initialized with the necessary Z -reference, and *passive* objects without it, and that messages in transit initially contain the reference also. To maintain R1 and R2 during the distributed computation, each transmission of a message copies the Z -reference, and processes delete

\mathbf{S}_p : { $state_p = active$ }
 send a message $\langle M, \langle \mathbf{cop} Z \rangle \rangle$

\mathbf{R}_p : { A basic message has arrived }
 receive message $\langle M, \langle \mathbf{cop} Z \rangle \rangle$; $state_p := active$;
 insert the Z -reference

\mathbf{I}_p : { $state_p = active$ }
 $state_p := passive$;
 delete all Z -references

Algorithm 3: The basic actions augmented with reference manipulation.

their Z -references when they become *passive*. More explicitly, the actions to be carried out by A_p are modified as presented in Algorithm 3.

With these modifications R1 and R2 are maintained indeed. R1 is maintained because Z -references are deleted in action \mathbf{I}_p , and inserted in action \mathbf{R}_p . The latter is possible because the message contains a Z -reference by R2. R2 is maintained because in action \mathbf{S}_p a Z -reference is included in every message. This is possible because only *active* objects send messages, and these objects contain a Z -reference by R1. Thus R1 and R2 are maintained during the computation, and by Theorem 2.1 IND holds. To arrive at a termination detection algorithm, superimpose a garbage collection algorithm upon the objects as described. The garbage collection algorithm is then modified so as to inform the objects A_p when it identifies Z as garbage. (When receiving this notice, the root objects enter the *terminated* state. However, we omit this trivial operation from the description of the algorithms that will follow.)

Theorem 2.2 *The algorithm as constructed satisfies conditions T1 and T2.*

Proof. To prove T1, assume any process enters the *terminated* state. This happens upon notice that Z is collected. By the correctness of the garbage collection algorithm (condition G1) this implies that Z is garbage. By IND the termination condition holds.

To prove T2, assume the termination condition holds. By IND, Z is garbage, hence, by the liveness of the garbage collector (condition G2) Z will eventually be collected. Notice of this will be sent to the processes, and these will enter the *terminated* state within finite time. \square

It was remarked in Section 1.2.2 that garbage collectors of the reference counting type are not able to collect cyclic structures of garbage, which may possibly harm the liveness of the termination detection algorithm. It is, however, easily seen that Z is not part of such a cyclic structure, and in fact the

following, stronger equivalence holds.

There are no references to $Z \Leftrightarrow$ the termination condition holds.

Hence, Theorem 2.2 also applies if the garbage collection algorithm satisfies only G1 and G3 (as is the case for a reference counting algorithm).

Summary of the Transformation. The construction of a termination detection algorithm is summarized in the following four steps.

1. Form the set \mathcal{O} of objects, consisting of the root objects A_p and one indicator object Z .
2. Superimpose upon the actions of the basic computation the handling of the Z -reference.
3. Superimpose upon this combined algorithm a garbage collection algorithm.
4. Replace the collection of Z by a notification of termination.

3 Examples of the Transformation

The transformation described in the previous section can in principle be applied to any garbage collection scheme, of the reference counting as well as the mark-and-sweep type, or working according to other principles. In the next three subsections “simple” distributed, weighted, and local reference counting are considered and corresponding termination detection algorithms are derived. In Section 3.4 the transformation is applied to a mark-and-sweep garbage collector.

3.1 Distributed Reference Counting

In this section we show how the distributed reference counting algorithm of Lermen and Maurer [16] can be transformed into a termination detection algorithm. For each non-root object o a reference count RC_o is maintained. When an o -reference is copied from object p to object q , p sends to q a copy message $\langle \mathbf{cop} \ o \rangle$ and to o an increment message $\langle \mathbf{inc} \ o, q \rangle$. When an o -reference is deleted by p , p sends a delete message (or *decrement* message) $\langle \mathbf{dec} \ o \rangle$ to o .

The scheme is complicated because of the possibility that $\langle \mathbf{dec} \ o \rangle$ and $\langle \mathbf{inc} \ o, q \rangle$ messages may arrive at o in a different order than they are sent. This is possible even in a system where message communication is FIFO, because these messages may be sent from different objects. (It is not possible, however, if message communication is synchronous or causally ordered [6].) If an $\langle \mathbf{inc} \ o, q \rangle$ message is overtaken by a $\langle \mathbf{dec} \ o \rangle$ message, RC_o may temporarily drop to 0, causing o to be collected while it is reachable. It is possible to overcome this problem by first sending the $\langle \mathbf{inc} \ o, q \rangle$ message and blocking the

sender until it receives an acknowledgement from o (and only then send the copy message). Interestingly, Lermen and Maurer found a solution which avoids the delay caused by the acknowledgement.

3.1.1 Description of the Scheme

The necessary synchronization between $\langle \mathbf{dec} \ o \rangle$ and $\langle \mathbf{inc} \ o, q \rangle$ messages is achieved by a two-way strategy. First, object o learns about the creation of a certain o -reference *before* it learns about the deletion of this o -reference. Second, o learns about the creation of all *copies* of a certain o -reference before it learns about the deletion of this o -reference. Lermen and Maurer show that this “causality preserving” protocol indeed implies the safety of the scheme (as stated in Theorem 3.1 further down).

To implement the first part, o sends to q an acknowledgement $\langle \mathbf{ack} \ o \rangle$ for the message $\langle \mathbf{inc} \ o, q \rangle$ it receives from p when an o -reference is copied from p to q . Note that the communication scheme is triangular: p sends $\langle \mathbf{cop} \ o \rangle$ to q and $\langle \mathbf{inc} \ o, q \rangle$ to o , o sends $\langle \mathbf{ack} \ o \rangle$ to q , and q receives both a $\langle \mathbf{cop} \ o \rangle$ message (from p) and an $\langle \mathbf{ack} \ o \rangle$ message (from o). The $\langle \mathbf{ack} \ o \rangle$ message informs q that o has learned about the creation of its o -reference. q deletes an o -reference only if it is an acknowledged reference, that is, q has received an $\langle \mathbf{ack} \ o \rangle$ message for it. q maintains a count both of the number of its acknowledged references to o and of the number of “surplus” copy messages it has received.

To implement the second part, Lermen and Maurer assume a FIFO discipline on the links. As p sends $\langle \mathbf{inc} \ o, q \rangle$ messages concerning copies of its o -reference earlier than the $\langle \mathbf{dec} \ o \rangle$ message, this ensures indeed that o is informed about the creation of copies before it learns about the deletion of the reference.

A formal description of the actions in the scheme is given as Algorithm 4. The actions to create a new object are omitted, because creation of objects does not occur when the scheme is used for termination detection. Each object p keeps, besides its multiset of references, for each non-root object o the two variables $aR_p(o)$, the number of acknowledged o -references, and $iR_p(o)$, the difference between the number of $\langle \mathbf{cop} \ o \rangle$ messages and the number of $\langle \mathbf{ack} \ o \rangle$ messages received by p . The conditional assignments in actions \mathbf{RC}_p and \mathbf{RA}_p cause p to increase its count of acknowledged o -references exactly in one of the following cases:

1. A $\langle \mathbf{cop} \ o \rangle$ message is received while there are unmatched $\langle \mathbf{ack} \ o \rangle$ messages ($iR_p(o) < 0$).
2. An $\langle \mathbf{ack} \ o \rangle$ message is received while there are unacknowledged o -references ($iR_p(o) > 0$).

Under the rules for the computation of the objects, an object already holding an o -reference may receive yet another $\langle \mathbf{cop} \ o \rangle$ message. In order to send enough $\langle \mathbf{dec} \ o \rangle$ messages in such a case, the action \mathbf{DR}_p is executed sufficiently

-
- CR_p**: { p is reachable and holds an o -reference }
 send $\langle \mathbf{cop} \ o \rangle$ to q ; send $\langle \mathbf{inc} \ o, q \rangle$ to o
- RC_p**: { A $\langle \mathbf{cop} \ o \rangle$ message has arrived at p }
 receive $\langle \mathbf{cop} \ o \rangle$; insert the o -reference ;
if $iR_p(o) < 0$ **then** $aR_p(o) := aR_p(o) + 1$;
 $iR_p(o) := iR_p(o) + 1$
- DR_p**: { $aR_p(o) > 0$ }
 send $\langle \mathbf{dec} \ o \rangle$ to o ; $aR_p(o) := aR_p(o) \Leftrightarrow 1$;
 delete an o -reference
- RI_o**: { An $\langle \mathbf{inc} \ o, q \rangle$ message has arrived at o }
 receive $\langle \mathbf{inc} \ o, q \rangle$;
 $RC_o := RC_o + 1$; send $\langle \mathbf{ack} \ o \rangle$ to q
- RD_o**: { A $\langle \mathbf{dec} \ o \rangle$ message has arrived at o }
 receive $\langle \mathbf{dec} \ o \rangle$; $RC_o := RC_o \Leftrightarrow 1$;
if $RC_o = 0$ **then** collect o
- RA_p**: { An $\langle \mathbf{ack} \ o \rangle$ message has arrived at p }
 receive $\langle \mathbf{ack} \ o \rangle$;
if $iR_p(o) > 0$ **then** $aR_p(o) := aR_p(o) + 1$;
 $iR_p(o) := iR_p(o) \Leftrightarrow 1$

Algorithm 4: The Lermen–Maurer scheme.

often to make $aR_p(o)$ equal to 0. An $\langle \mathbf{ack} \ o \rangle$ message may even arrive when the reference is no longer needed, in which case the execution of **DR_p** is deferred until the receipt of the $\langle \mathbf{ack} \ o \rangle$ message (if **RA_p** results in $aR_p(o) > 0$). The termination detection algorithm, obtained from this scheme in the next subsection, must also allow multiple execution of the corresponding action, and may defer it until an $\langle \mathbf{ack} \ o \rangle$ message has been received (see also the remarks at the end of Section 3.1.2).

Initially there are only acknowledged references ($iR_p(o) = 0$ for all p, o), the reference counts correctly reflect their number ($RC_o = \sum_p aR_p(o)$), and no messages are in transit. The correctness of the protocol, as expressed in the following theorem, was proved by Lermen and Maurer in [16].

Theorem 3.1 *Algorithm 4 is a correct reference counting garbage collection algorithm, that is, it satisfies G1 and G3.*

3.1.2 Transformation into a Termination Detection Algorithm

In this section a termination detection algorithm is derived from the Lermen–Maurer scheme. A discussion of the properties of the derived algorithm, called the *Activity Counting* algorithm, is deferred to Section 3.1.3. The termination detection algorithm is derived in the four steps described at the end of Section 2.

1. The set \mathcal{O} of objects consists of the objects A_p and the indicator object Z .
2. Superimpose upon the actions of the basic computation the handling of the Z -reference. Algorithm 3 is obtained.
3. Superimpose upon Algorithm 3 the garbage collection algorithm of Lermen and Maurer. To this end, the \mathbf{CR}_p action is included in the \mathbf{S}_p action, the \mathbf{RC}_p action is included in the \mathbf{R}_p action, and the \mathbf{DR}_p action is included in the \mathbf{I}_p action. In all cases Z is substituted for o . This operation yields Algorithm 5.
4. Replace the collection of Z by a notification of termination. This is done by substituting

send $\langle \mathbf{term} \rangle$ to all A_p

for “collect Z ” in action \mathbf{RD}_Z . Upon receipt of this message, the processes enter the *terminated* state.

The derivation of the termination detection algorithm is now complete. Finally, preserving its correctness, the algorithm can be simplified. Because there is only one non-root object, the subscript Z may be dropped from all variables. Furthermore, the handling of the Z -reference only serves to lead the garbage collection scheme in its actions. Now that these actions have been correctly connected to the actions of the basic distributed computation, this reference handling can be removed. The resulting Activity Counting algorithm is given in Algorithm 6. The initial conditions for this algorithm are: $iR_p = 0$ for all p ; $aR_p = 0$ if p is *passive*, and $aR_p > 0$ if p is *active*; $RC = \sum_p aR_p$; and no messages are in transit.

The correctness of the algorithm, as expressed in the following theorem, is a direct consequence of the properties of the Lermen–Maurer scheme and the validity of our transformation.

Theorem 3.2 *The Activity Counting algorithm is a correct termination detection algorithm.*

Proof. According to Theorem 3.1, the Lermen–Maurer scheme satisfies G1 and G3, and hence by Theorem 2.2 (and the remark following it), the derived algorithm satisfies T1 and T2. \square

S_p : $\{ state_p = active \}$ (* This implies p holds a Z -reference *)
 send a message $\langle M, \langle \mathbf{cop} Z \rangle \rangle$ to q ; send $\langle \mathbf{inc} Z, q \rangle$ to Z

R_p : $\{ \text{A basic message has arrived} \}$
 receive message $\langle M, \langle \mathbf{cop} Z \rangle \rangle$; $state_p := active$;
 insert the Z -reference ;
if $iR_p(Z) < 0$ **then** $aR_p(Z) := aR_p(Z) + 1$;
 $iR_p(Z) := iR_p(Z) + 1$

I_p : $\{ aR_p(Z) > 0 \}$
 $state_p := passive$;
 send $\langle \mathbf{dec} Z \rangle$ to Z ; $aR_p(Z) := aR_p(Z) \Leftrightarrow 1$;
 delete a Z -reference from the references of A_p

RI_Z : $\{ \text{An } \langle \mathbf{inc} Z, q \rangle \text{ message has arrived at } Z \}$
 receive $\langle \mathbf{inc} Z, q \rangle$;
 $RC_Z := RC_Z + 1$; send $\langle \mathbf{ack} Z \rangle$ to q

RD_Z : $\{ \text{A } \langle \mathbf{dec} Z \rangle \text{ message has arrived at } Z \}$
 receive $\langle \mathbf{dec} Z \rangle$; $RC_Z := RC_Z \Leftrightarrow 1$;
if $RC_Z = 0$ **then** collect Z

RA_p : $\{ \text{An } \langle \mathbf{ack} Z \rangle \text{ message has arrived at } p \}$
 receive $\langle \mathbf{ack} Z \rangle$;
if $iR_p(Z) > 0$ **then** $aR_p(Z) := aR_p(Z) + 1$;
 $iR_p(Z) := iR_p(Z) \Leftrightarrow 1$

Algorithm 5: The basic computation with the Lermen–Maurer scheme.

A process may receive an activation message when it is already *active*, in which case it does not become “even more *active*”. In order to send enough $\langle \mathbf{dec} \rangle$ messages, the process must later execute action I_p as many times as it has received activation messages. This is taken care of by the guard of this action, based on the proper administration of the $\langle \mathbf{ack} \rangle$ and $\langle M \rangle$ messages received (variables aR_p and iR_p). When the process is *passive*, but has positive aR_p , it eventually executes I_p and sends a $\langle \mathbf{dec} \rangle$ message.

3.1.3 Discussion of the Algorithm

The principle of the Activity Counting algorithm is simple. When a process activates another process, it informs the central controller Z by sending an increment message $\langle \mathbf{inc} q \rangle$, and when a process becomes *passive*, it informs Z by sending a decrement message $\langle \mathbf{dec} \rangle$. The controller tries to keep an account

S_p: { *state_p* = *active* }
 send a message $\langle M \rangle$ to q ; send $\langle \mathbf{inc} \ q \rangle$ to Z

R_p: { A basic message has arrived }
 receive message $\langle M \rangle$; *state_p* := *active* ;
 if $iR_p < 0$ **then** $aR_p := aR_p + 1$;
 $iR_p := iR_p + 1$

I_p: { $aR_p > 0$ }
 state_p := *passive* ;
 send $\langle \mathbf{dec} \rangle$ to Z ; $aR_p := aR_p \ominus 1$

RI_Z: { An $\langle \mathbf{inc} \ q \rangle$ message has arrived at Z }
 receive $\langle \mathbf{inc} \ q \rangle$;
 $RC := RC + 1$; send $\langle \mathbf{ack} \rangle$ to q

RD_Z: { A $\langle \mathbf{dec} \rangle$ message has arrived at Z }
 receive $\langle \mathbf{dec} \rangle$; $RC := RC \ominus 1$;
 if $RC = 0$ **then** send $\langle \mathbf{term} \rangle$ to all A_p

RA_p: { An $\langle \mathbf{ack} \rangle$ message has arrived at p }
 receive $\langle \mathbf{ack} \rangle$;
 if $iR_p > 0$ **then** $aR_p := aR_p + 1$;
 $iR_p := iR_p \ominus 1$

Algorithm 6: The Activity Counting termination detection algorithm.

of the number of “currently” *active* processes by counting the increment and decrement messages. When enough decrement messages have been received to balance the increment messages and the initially *active* processes, it signals termination. Unfortunately, the possible delay between a basic action and its registration by Z render this over-simplified scheme incorrect as the following example shows.

1. Assume only p is *active* and $RC = 1$.
2. p sends an activation message to q and an increment message to Z .
3. q receives the activation message and becomes *active*. Then q becomes *passive* again and sends a decrement message to Z .
4. Z receives the decrement message (before the increment message), and RC drops to 0 while p is still *active*.

The Activity Counting algorithm takes care of this and similar scenarios, because the sending of the decrement message is deferred until an acknowledgement message $\langle \mathbf{ack} \rangle$ has been received. This implies that Z receives the decrement message *after* it has received the corresponding increment message.

As the Activity Counting algorithm was derived from the Lermen–Maurer reference counting scheme, it inherits properties of the latter algorithm, which will now be discussed.

1. Message Complexity. The message overhead of the new algorithm is considerable: for each basic message of the computation the algorithm adds up to *three* control messages ($\langle \mathbf{inc} \ q \rangle$, $\langle \mathbf{ack} \rangle$ and $\langle \mathbf{dec} \rangle$). A worst case lower bound for this overhead of *one* control message per basic message was proved by Chandy and Misra [4], and this bound is achieved by the algorithm of Dijkstra and Scholten [11].

2. FIFO Discipline on Links. For the correctness of the algorithm it is required that links deliver messages in the order they were sent.

3. Central Controller. The central object Z acts as a central controller in the Activity Counting algorithm. For *each* single transmission of a basic message up to *two* actions are necessary in Z (\mathbf{RI}_Z and \mathbf{RD}_Z). The central process may become a bottleneck in the computation and slow down the operation of the entire system. The local reference counting algorithm described later in Section 3.3 also relies on a central controller, but does not require its cooperation with each basic message transmission.

4. Inhibition. A control algorithm is said to be *inhibitory* if it may temporarily disable actions of the basic computation. In the Lermen–Maurer scheme this is the case for the delete action, which is deferred until the object has an acknowledged version of the reference. As a consequence, in the Activity Counting algorithm becoming *passive* is only allowed if $aR_p(o) > 0$, thus formally the algorithm is inhibitory. This disadvantage can be overcome by a slight modification of the algorithm as follows. The object may delete any reference, but the $\langle \mathbf{dec} \ o \rangle$ message is held back if the acknowledgement for the reference was not yet received. A similar modification makes the Activity Counting algorithm non-inhibitory.

3.1.4 Related Algorithms

The Vector Counting Algorithm. Through the use of $\langle \mathbf{ack} \rangle$ messages (and FIFO channels) the Activity Counting algorithm guarantees that Z always has a causally consistent (though possibly slightly outdated) view of the number of active (or activated) processes in the system. By keeping more information in Z , however, it is also possible to achieve this without actually sending $\langle \mathbf{ack} \rangle$

messages. For this purpose Z keeps a vector (i.e., an integer array) V with one component for each process. Whenever Z receives an $\langle \mathbf{inc} \ q \rangle$ message it increments q 's component of V instead of sending an $\langle \mathbf{ack} \rangle$ message: $V[q] := V[q] + 1$. Action \mathbf{RA}_p and the variables iR_p are no longer necessary: a process increments aR_p when it receives a basic message and sends aR_p $\langle \mathbf{dec} \rangle$ messages to Z immediately when it becomes passive. When Z receives a $\langle \mathbf{dec} \rangle$ message from q it decrements the corresponding component of V : $V[q] := V[q] \ominus 1$.

Notice that temporarily $V[q]$ might become negative—this is the case if Z receives a $\langle \mathbf{dec} \rangle$ message *before* it receives the corresponding $\langle \mathbf{inc} \ q \rangle$ message. This is precisely the situation which is avoided by use of $\langle \mathbf{ack} \rangle$ messages in the original Activity Counting algorithm. Because Z 's view is inconsistent when a component of V is negative, nothing is deduced in that case. It follows that Z can signal termination when V becomes the null vector, and the RC counter is no longer necessary. Some further optimizations (e.g., batching $\langle \mathbf{inc} \ q \rangle$ and $\langle \mathbf{dec} \rangle$ messages) yield a centralized variant of the so-called *Vector Counter* termination detection algorithm proposed in [17]. This algorithm has lower message overhead than the Activity Counting algorithm, does not rely on the FIFO property, and can easily be realized in a distributed way as well.

Variations of the Lermen–Maurer Scheme. Two variants of the Lermen–Maurer scheme were proposed by Rudalics [25]. We describe informally his *three message protocol*, which does not rely on FIFO links. In the Lermen–Maurer scheme object q must receive an $\langle \mathbf{ack} \ o \rangle$ acknowledgement when an o -reference has been copied to it, but in the three message protocol an object p must receive an acknowledgement when it has initiated a copy of an o -reference. A delete message for the o -reference may be sent only when all acknowledgements have been received, and to this end an acknowledgement counter is added to each reference. The protocol works as follows.

1. To copy an o -reference to q , p increments the acknowledgement count of its reference and sends an increment message to o .
2. On receipt of this message, o increments its reference count and sends a copy message to q .
3. On receipt of the copy message, q inserts the reference and sends an acknowledgement to p .
4. On receipt of this acknowledgement, p decrements the acknowledgement count of the o -reference.

When p deletes the reference, it holds back the delete message until all acknowledgements have been received. To this end, the references are all installed with acknowledgement count equal to 1, and deletion of the reference is done by decrementing the count. When the count drops to zero (either by deletion of the reference or by receipt of an acknowledgement) the delete message is sent.

A drawback of this algorithm is that the new reference can only be installed after two messages have been propagated (one from p to o and one from o to q). In Rudalics’ *four message protocol* p also sends a copy message to q directly, and q installs the reference when it receives a copy message either from p or from o . The acknowledgement is sent when both copy messages have been received.

In a similar way as for the Lermen–Maurer scheme a termination detection algorithm can be derived from the three and four message protocols by the transformation of Section 2.

3.2 Weighted Reference Counting

In this section we consider the transformation of a garbage collection algorithm based on *weighted reference counting* (WRC). The resulting termination detection algorithm turns out to be an already known algorithm: it was proposed by Mattern [18].

As mentioned in the introduction of Section 3.1, $\langle \mathbf{dec} \ o \rangle$ and $\langle \mathbf{inc} \ o, q \rangle$ messages in the “simple” distributed reference counting scheme must be synchronized because their reordering may render the scheme unsafe. This need for synchronization can be avoided using weighted reference counting. In this variant each reference has a positive *weight*. The reference count of an object o now represents the total weight of all existing o –references rather than their number. (We continue to use the word reference *count* although it may no longer be completely appropriate.) When a reference is copied, its weight is *split* among the existing and the new reference. Thus, although the *number* of references increases, the *weight* remains the same, and the reference count need not be incremented and no $\langle \mathbf{inc} \ o, q \rangle$ message need be sent. When an object deletes an o –reference, a decrement message is sent to o , returning the weight of the deleted reference. Upon receipt of such a message, o subtracts the weight from its reference count. The reference count monotonically decreases, and the order in which control messages (i.e., delete messages) arrive at the object becomes irrelevant. The object can be collected when its reference count drops to zero.

3.2.1 Description of the Scheme

Distributed weighted reference counting schemes have been given by Bevan [2], Watson and Watson [36], and others. The principle was attributed to Weng [35]. In its description, see Algorithm 7, again the mechanism to create new objects is omitted. An o –reference is now a tuple (o, w) , where w denotes the weight of the reference. Initially for each non–root object o , the reference count RC_o equals the sum of the weights of all existing o –references.

Theorem 3.3 *Algorithm 7 is a correct reference counting garbage collection algorithm, that is, it satisfies G1 and G3.*

Proof. A correctness proof and analysis of the scheme is given by Bevan [2] and by Watson and Watson [36] and is based on invariance of the following two assertions:

CR_p: { p is reachable and holds a reference (o, w) }
 send $\langle \mathbf{cop} \ o, w/2 \rangle$ to q ; $w := w/2$

RC_p: { A message $\langle \mathbf{cop} \ o, w \rangle$ has arrived }
 receive $\langle \mathbf{cop} \ o, w \rangle$;
if p has an o -reference
 then add w to its weight
 else insert the o -reference with weight w

DR_p: { p holds a reference (o, w) }
 send $\langle \mathbf{dec} \ o, w \rangle$ to o ; delete the o -reference

RD_o: { A $\langle \mathbf{dec} \ o, w \rangle$ message has arrived at o }
 receive $\langle \mathbf{dec} \ o, w \rangle$; $RC_o := RC_o \oplus w$;
if $RC_o = 0$ **then** collect o

Algorithm 7: The weighted reference counting scheme.

1. Each reference has a positive weight; each delete message contains a positive weight.
2. $RC_o = \sum_{R=(o,w)} w + \sum_{D=(\mathbf{dec} \ o,w)} w$, where R ranges over all o -references in existence (including those in copy messages) and D ranges over all delete messages in transit.

□

3.2.2 Transformation into a Termination Detection Algorithm

To transform the garbage collection scheme into a termination detection algorithm we apply the four step construction of Section 2. Steps 1, 2, and 4 are as in Section 3.1.2. In step 3 the actions of the weighted reference counting scheme are superimposed on the program resulting from step 2 (Algorithm 3). To this end, action **CR_p** is included in action **S_p**, action **RC_p** is included in action **R_p**, and action **DR_p** is included in action **I_p**. Again for o the object Z is substituted. This results in Algorithm 8.

The same simplifications as in Section 3.1.2 can be made. The actual handling of the Z reference can be removed, and instead we simply equip every process p with a variable W_p , representing the weight of p 's (virtual) Z -reference (0 if p has no such reference). The subscript Z is dropped. This finally results in Algorithm 9, which is known as the *Credit Recovery algorithm* [18]. The initial conditions for this algorithm are: $W_p = 0$ if p is *passive*; $W_p > 0$ if p is *active*; $RC = \sum_p W_p$; and no messages are in transit.

S_p: { $state_p = active$ } (* Thus p has a Z -reference (Z, w) *)
 send a message $\langle M, \langle \mathbf{cop} Z, w/2 \rangle \rangle$ to q ; $w := w/2$

R_p: { A basic message has arrived }
 receive message $\langle M, \langle \mathbf{cop} Z, w \rangle \rangle$; $state_p := active$;
 if p has a Z -reference
 then add w to its weight
 else insert the Z -reference with weight w

I_p: { $state_p = active$ }
 $state_p := passive$;
 send $\langle \mathbf{dec} Z, w \rangle$ to Z ; delete the Z -reference

RD_Z: { A $\langle \mathbf{dec} Z, w \rangle$ message has arrived at Z }
 receive $\langle \mathbf{dec} Z, w \rangle$; $RC_Z := RC_Z \ominus w$;
 if $RC_Z = 0$ **then** collect Z

Algorithm 8: The basic computation with weighted reference counting.

S_p: { $state_p = active$ } (* Thus $W_p > 0$ *)
 send a message $\langle M, W_p/2 \rangle$ to q ; $W_p := W_p/2$

R_p: { A basic message has arrived }
 receive message $\langle M, W \rangle$; $state_p := active$;
 $W_p := W_p + W$

I_p: { $state_p = active$ }
 $state_p := passive$;
 send $\langle \mathbf{dec} W_p \rangle$ to Z ; $W_p := 0$

RD: { A $\langle \mathbf{dec} W \rangle$ message has arrived at Z }
 receive $\langle \mathbf{dec} W \rangle$; $RC := RC \ominus W$;
 if $RC = 0$ **then** send $\langle \mathbf{term} \rangle$ to all A_p

Algorithm 9: The Credit Recovery termination detection algorithm.

Theorem 3.4 *The Credit Recovery algorithm is a correct termination detection algorithm.*

Proof. According to Theorem 3.3, the weighted reference counting scheme satisfies G1 and G3, and hence by Theorem 2.2 (and the remark following it),

the derived algorithm satisfies T1 and T2. \square

3.2.3 Discussion of the Algorithm

Alternative Implementations of Action \mathbf{RC}_p . Action \mathbf{RC}_p of Algorithm 7 adds the weight of the received reference if p already has an o -reference. There are two alternative implementations of the algorithm. First, p may return the received weight to o immediately in a $\langle \mathbf{dec} \ o, w \rangle$ message. Second, p may store the weight separately and thus keep a non-empty *set* of weights, one for each o -reference, rather than a single weight. Both alternatives maintain the two invariants of the weighted reference counting algorithm and are therefore also applicable to the credit recovery termination detection algorithm. A consequence of the two alternative strategies is that all weights in the system are always (negative) powers of 2, and can thus be represented concisely by their negative logarithm.

Weight Underflow. The implementation of the weighted reference counting scheme faces a difficulty that has not yet been discussed, and it is not surprising that the Credit Recovery algorithm faces a similar difficulty. The problem arises because weights are represented in a finite number of bits: thus there is a smallest possible positive weight, and if a reference of this weight is copied its weight cannot be split. The problem in the Credit Recovery algorithm arises when a process with the smallest possible positive value of W_p sends a message.

Also the solutions to these difficulties are similar in the two algorithms. In the weighted reference counting algorithms, a new *indirection object* is created with a *maximal* reference count, and the original reference is replaced by a reference to the indirection object, with maximal weight. Next it can be copied without difficulties. In the Credit Recovery algorithm, a process negotiates with Z to exchange its (minimal) credit for a new, maximal credit. Then it can send the message. The operation results in an increase in the reference count of Z .

These additions to the algorithms do not make them inefficient or impractical, because in both algorithms weight underflow is supposed to be a very rare event. As remarked above, it can be arranged that all weights are powers of 2, and can be represented by their (negative) logarithm. When copying a reference, $W := W + 1$ is executed instead of $W := W/2$, and the probability of overflow in W should not be much greater than the overflow probability in classical reference counting schemes.

Generational Reference Counting. Another reference counting principle, called *generational reference counting*, has been proposed by Goldberg [14]. This scheme, like weighted reference counting, avoids $\langle \mathbf{inc} \ o, q \rangle$ messages, but uses a different strategy for this. The communication pattern is the same as for the weighted reference counting scheme: a control message is sent to o only upon deletion of an o -reference.

\mathbf{S}_p : { $state_p = active$ }
 send $(M, gen_p + 1)$ to q ;
 $sons_p := sons_p + 1$

\mathbf{R}_p : { A basic message (M, g) has arrived }
if $state_p = active$
 then send $\langle \mathbf{pas} \ g, 0 \rangle$ to Z
 else $state_p, gen_p, sons_p := active, g, 0$

\mathbf{I}_p : { $state_p = active$ }
 send $\langle \mathbf{pas} \ gen_p, sons_p \rangle$ to Z ;
 $state_p := passive$

\mathbf{RP} : { A $\langle \mathbf{pas} \ g, s \rangle$ message has arrived at Z }
 receive $\langle \mathbf{pas} \ g, s \rangle$; $ActCount[g] := ActCount[g] \ominus 1$;
 $ActCount[g + 1] := ActCount[g + 1] + s$;
if $\forall i : ActCount[i] = 0$ **then** send $\langle \mathbf{term} \rangle$ to all A_p

Algorithm 10: The Generational termination detection algorithm.

Each reference has a *generation* number, where a copy of a reference with generation number i has generation number $i + 1$. An object keeps separate reference counts for each generation. When an o -reference is deleted, a $\langle \mathbf{dec} \ o, g, s \rangle$ message is sent to o , reporting the generation number g of the deleted reference and the number s of copies (with generation number $g + 1$) that were made of the reference. We omit a full description of the scheme and the (straightforward) transformation into a termination detection algorithm.

The resulting *Generational termination detection algorithm* is given in Algorithm 10. Each *active* process p has a *generation number* gen_p and a counter $sons_p$ to count how many activation messages it has sent (in its current *active* period). When a process becomes *passive* it sends to Z a $\langle \mathbf{pas} \ gen_p, sons_p \rangle$ message to inform it that it has become *passive* and sent $sons_p$ activation messages (of generation $gen_p + 1$).

The central controller Z maintains an array $ActCount$ of integers, where $ActCount[g]$ counts the *active* processes of the g^{th} generation. When a $\langle \mathbf{pas} \ g, s \rangle$ message is received, $ActCount[g]$ is decremented and $ActCount[g + 1]$ is increased by s , the number of newly reported activations of generation $g + 1$. Initially a process p is either *active* with $gen_p = 1$ or *passive*, $ActCount[1]$ equals the number of *active* processes, $ActCount[i] = 0$ for $i > 1$, and no messages are underway.

Theorem 3.5 *The Generational termination detection algorithm is a correct termination detection algorithm.*

Proof. This result follows from the correctness of the generational reference counting scheme (as demonstrated by Goldberg [14]) and Theorem 2.2. \square

We present this algorithm as another illustration of our transformation, but Goldberg’s remark: “it is not clear if there is any advantage to using the generational reference counting scheme instead of the weighted reference counting scheme” seems to apply equally to the resulting termination detection algorithms.

3.3 Local Reference Counting

In [15] Ichisugi and Yonezawa present an interesting distributed garbage collection scheme they call *local reference counting* (LRC). Basically the same idea was found independently by Piquer [24] (“indirect reference counting”) and by Rudalics [26]. The proposed scheme assumes that objects are partitioned into groups and each group is mapped onto a distinct computing node. Thus, the model used here is slightly different from the model described at the beginning of Section 1.2. In particular, a difference between “local” and “remote” objects and references is now made. An interesting feature of the LRC garbage collection scheme is that it supports the migration of objects from one node to another node very easily. The interested reader is referred to [24].

3.3.1 Description of the Scheme

For each (local or remote) object o which is referenced by a local object of node N , N has a *local reference counter* $LRC_N(o)$. This counter is incremented when an o -reference is locally or remotely copied, and it is decremented when a local o -reference is deleted or when all offsprings of a remotely copied o -reference have been deleted. On a local variable $FIRST_N(o)$, N keeps the identity of the node from which it first received an o -reference. $FIRST_N(o)$ is undefined (\emptyset) if o was created by an object residing on N . If $LRC_N(o)$ drops to 0, o is garbage if it was created by a local object. Otherwise the node $FIRST_N(o)$ is informed by a $\langle \mathbf{dec} \rangle$ message that N is unaware of any remaining o -references.

For our purpose it is appropriate to assume that exactly one object resides on each node (which means that we identify nodes and objects). The set of rules for garbage identification can then be stated in a slightly adapted form (compared to the rules given in [15]) as in Algorithm 11. Obviously, action \mathbf{DZ}_p can be appended to actions \mathbf{DR}_p and \mathbf{RD}_p guarded by a test “**if** $LRC_p(o) = 0$ **then** ...”. Notice that when an object that already has an o -reference receives another o -reference (action \mathbf{RC}_p) a $\langle \mathbf{dec} \ o \rangle$ message is returned to the sender although the reference is installed, and the local reference counter is incremented.

Compared to the Lermen-Maurer scheme (Section 3.1), the LRC algorithm has the important advantage that no extra messages besides the copy message (or decrement message) are necessary when a reference is copied (or deleted). Furthermore, the FIFO property is not required. Compared to the weighted reference counting scheme (Section 3.2), the handling of the local reference

CR_p: { p holds an o -reference }
 send $\langle \mathbf{cop} \ o \rangle$ to q ; $LRC_p(o) := LRC_p(o) + 1$

RC_p: { A $\langle \mathbf{cop} \ o \rangle$ message sent by q has arrived at p }
 receive $\langle \mathbf{cop} \ o \rangle$; insert the o -reference ;
if $LRC_p(o) = 0$
 then $LRC_p(o) := 1$; $FIRST_p(o) := q$;
 else send $\langle \mathbf{dec} \ o \rangle$ to q ; $LRC_p(o) := LRC_p(o) + 1$

DR_p: { p holds an o -reference }
 delete the o -reference ; $LRC_p(o) := LRC_p(o) \ominus 1$

RD_p: { A $\langle \mathbf{dec} \ o \rangle$ message has arrived }
 receive $\langle \mathbf{dec} \ o \rangle$; $LRC_p(o) := LRC_p(o) \ominus 1$

DZ_p: { $LRC_p(o)$ has just dropped from 1 to 0 }
if $FIRST_p(o) \neq \emptyset$
 then send $\langle \mathbf{dec} \ o \rangle$ to $FIRST_p(o)$;
 else collect o

Algorithm 11: The local reference counting scheme.

counters is somewhat simpler than the accumulation of arbitrary small fragments on the weight reference counter RC_o . A drawback of the LRC scheme is that a node (or an object) usually has many counters (whereas in the weighted reference counting scheme a single counter per object is sufficient), possibly one for each reference that ever existed on that node.

It is also interesting to compare the LRC scheme to the generational reference counting principle because the local reference counters can also be regarded as some sort of generational counters. These counters, however, are not kept at the referenced object o but are distributed in a tree-like manner over all nodes (or objects) which have an o -reference. The levels of this tree correspond to the generation number.

3.3.2 Transformation into a Termination Detection Algorithm

For the transformation of the LRC garbage collection algorithm into a termination detection algorithm we proceed as before. Action **CR_p** is included in action **S_p**, action **RC_p** is included in action **R_p**, and action **DR_p** is included in action **I_p**. In order to avoid the keeping of more than one Z -reference per object (which could happen if an active process is reactivated), we delete a second

S_p: { *state_p* = *active* }
 send a message $\langle M \rangle$ to q ; $LRC_p := LRC_p + 1$

R_p: { A basic message $\langle M \rangle$ sent by q has arrived }
 receive message $\langle M \rangle$;
 if $LRC_p = 0$
 then $LRC_p := 1$; $FIRST_p := q$; *state_p* := *active* ;
 else send $\langle \mathbf{dec} \rangle$ to q ; **if** *state_p* = *passive*
 then *state_p* := *active* ; $LRC_p := LRC_p + 1$

I_p: { *state_p* = *active* }
 state_p := *passive* ;
 $LRC_p := LRC_p \Leftrightarrow 1$;
 if $LRC_p = 0$ **then** send $\langle \mathbf{dec} \rangle$ to $FIRST_p$

RD_p: { A $\langle \mathbf{dec} \rangle$ message has arrived }
 receive $\langle \mathbf{dec} \rangle$; $LRC_p := LRC_p \Leftrightarrow 1$;
 if $LRC_p = 0$ **then** send $\langle \mathbf{dec} \rangle$ to $FIRST_p$

Algorithm 12: The resulting LRC termination detection algorithm.

Z -reference immediately after it has been installed in action **RC_p**. Technically this means that the increment of the local reference counter is suppressed if an already active process executes action **RC_p** (or action **R_p** below).

We assume that the virtual object Z is created by the “environment” e before the start of the computation. Conceptually, the environment e can be regarded as a special process which is the only initially active process and which starts the computation by sending basic messages to other processes (the “initiators”). It does not receive basic messages and consequently its variable $FIRST_e$ is always undefined. Initially, no other process has a Z -reference. Therefore, $FIRST_p$ is never undefined for a regular process p if LRC_p eventually drops to 0. This simplifies action **DZ_p** which is appended to actions **RD_p** and **DR_p** (**I_p** resp.). As before, the handling of Z -references is omitted in the resulting Algorithm 12.

For the environment e , LRC_e must be initialized to 1 in order to reflect the existence of its virtual Z -reference before the start of the computation. Because $FIRST_e = \emptyset$, the last line of actions **I** and **RD** in the environment should be replaced by

if $LRC_e = 0$ **then** send $\langle \mathbf{term} \rangle$ to all A_p .

3.3.3 Discussion of the Algorithm

The resulting termination detection algorithm is basically identical to the termination detection scheme for diffusing computations presented by Dijkstra and Scholten [11]. The authors call the **<dec>** messages “signals” and keep two local counters C_p (sum of the *deficits* of p ’s incoming edges) and D_p (sum of the *deficits* of p ’s outgoing edges) for each process p . When a basic message is received C_p is incremented; C_p is decremented when a signal is sent. D_p is incremented when a basic message is sent (which is only allowed if $C_p > 0$) and decremented when a signal is received. According to Dijkstra and Scholten’s protocol, sending a signal is only possible if

$$\text{(SIG)} \quad C_p > 1 \text{ or } (C_p = 1 \text{ and } D_p = 0).$$

This means that if *before* the receipt of the message $C_p = 1$, p can immediately send a signal in reaction to an incoming message (and consequently reset C_p to 1). If this immediate reaction discipline is adopted, it is always the case that $C_p = 0$ or $C_p = 1$. Because $C_p = 0 \Rightarrow D_p = 0$ is a required invariant of the scheme, it follows from (SIG) that (besides the immediate reaction to an incoming message) the sending of a signal is constrained by the condition $C_p + D_p = 1$. After the sending of a signal the property $C_p + D_p = 0$ holds. It should be clear now that the local variable LRC_p of our derived algorithm plays the role of $C_p + D_p$; a signal is sent when LRC_p drops from 1 to 0. By (SIG) Dijkstra and Scholten’s scheme allows to delay the sending of a signal, but otherwise their scheme is identical to our derived algorithm. Since Dijkstra and Scholten prove the correctness of their scheme in [11], we directly have the following theorem.

Theorem 3.6 *The LRC termination detection algorithm is a correct termination detection algorithm.*

The LRC termination detection algorithm has interesting characteristics. For each basic message eventually one decrement control message is sent which means that the algorithm has optimal worst case message overhead [4]. The FIFO property of communication links is not required. In contrast to the previously derived Activity Counting algorithm (Section 3.1.2) and the Credit Recovery algorithm (Section 3.2.2) the possible bottleneck of a central controller Z (which in those schemes receives a control message each time a process becomes passive) is avoided. The memory overhead is small, each process only needs a pointer ($FIRST_p$) and a counter (LRC_p).

3.4 Mark-and-Sweep Garbage Collection

As explained in Section 1.2.2, mark-and-sweep garbage collectors operate in consecutive cycles. In each cycle first all reachable objects are marked, and subsequently all unmarked objects are reclaimed. In this section it is shown how the garbage collection algorithm of Ben-Ari [1] can be transformed into

a termination detection algorithm. (Actually, we use a variant of the algorithm described by Van de Snepscheut [29]). This algorithm was designed to run concurrently with a single processor mutating the references contained in memory cells. Thus the copying of a reference is a single atomic step, where we have assumed so far that it consists of the sending and receipt of a message. When using Ben–Ari’s collector, these two events must be assumed to be one single event. Therefore, in the remainder of this section we assume *synchronous* communication. Let n denote the number of processes.

3.4.1 Description of the Scheme

In each cycle initially all nodes are white, and the following is done in a cycle¹:

1. Color all roots gray.
2. Sequentially visit all nodes. For all gray nodes, color the nodes to which they have references gray.
3. Sequentially visit all nodes and count the number of gray nodes.
4. If more gray nodes were counted than in the previous round (more than the number of roots for the first round) go to step 2, otherwise to step 5.
5. Collect the white nodes and make all nodes white.

It is essential for the correctness of the algorithm that the basic program cooperates with the marker algorithm: whenever the basic program installs a new reference, it makes the object to which it points gray. No cooperation is required when the basic program deletes a reference. The correctness proof of this garbage collection scheme is quite involved; proofs were given by Ben–Ari in [1] and Van de Snepscheut [29].

3.4.2 Transformation into a Termination Detection Algorithm

In the scheme used to obtain a termination detection algorithm each object can have at most one reference, which is always a Z –reference. The transformation is straightforward. Rather than coloring the roots white at the end of each cycle and gray again at the beginning of the next one, assume that the roots are always gray by definition. The five steps of the algorithm are transformed as follows.

1. (Gray the roots.) The roots (i.e., the processes A_p) are always gray, so this step is skipped.
2. (Gray sons of gray nodes.) In this step the virtual object Z need not be visited as it has no sons. The processes A_p are visited by arranging the

¹We use the color “gray” instead of “black” (as in the original algorithm [1, 29]) because we need “black” for a different purpose further down.

processes in a (virtual) ring and passing a token along this ring. On this tour the token visits the processes in a “lazy” way: before actually visiting a process, it waits until the process is *passive*, and thus has no reference at all. This does not hinder the liveness of the termination detection, because there is no termination while a process is *active*. As a result of this strategy, no coloring is done in this round.

3. (Count gray nodes.) There are $n + 1$ processes, and the n roots are known to be gray. It only needs to be determined whether any root has grayed Z , which is the case if a Z -reference has been installed (by the basic program) since the beginning of this cycle. To this end, a process (i.e., a root) becomes *black* when it would gray Z according to the scheme (*viz.*, when becoming *active*). In order to see whether Z was grayed, the token again visits all processes, now testing whether any process is black.
4. (Cycle completed?) If the second tour of the token reveals that no process is black, then Z is white and the number of gray nodes is still n . In this case, go to step 5. If any node is black (i.e., Z is gray) there are now $n + 1$ non-white nodes and the original algorithm would jump to step 2 in order to (try to) gray more nodes. However, in our case this is useless, as no nodes are white anymore, and we decide to also terminate the cycle. As Z is gray, it cannot be collected, hence a new cycle of the collector must be started by resetting the color of all roots to gray and returning to step 2.
5. (Collect.) The collection phase is entered with n gray nodes, i.e., the virtual node Z is white, therefore termination can be signaled.

Essential for the termination detection algorithm is that a (gray) process becomes black when it activates another process. (The meaning of this blackening in the garbage collection scheme is “I grayed Z ”.) The termination detection algorithm repeatedly sends a token around the ring twice. In the first round the token only waits at each process until this process is *passive*. In the second round the token inspects the color of the processes and resets them to gray. If no process was black (i.e., Z is white), termination is concluded, otherwise a new double tour of the token is initiated.

3.4.3 Correctness of the New Algorithm

In this section we formally present the new termination detection scheme (see Algorithm 13), and prove its correctness by means of an invariant. Assume the processes are numbered from 0 through $n \Leftrightarrow 1$ and they have communication facilities so that process q can send control messages to process $q \Leftrightarrow 1 \pmod n$. The variable *tour* (values *first*, *second*) denotes whether the token is on the first or the second of the two tours, and t denotes the current position of the token. Processes have a color, stored in $color_q$ for process q . Instead of the colors black and gray used in the previous section, we prefer the colors black

T1: $\{ \text{tour} = \text{first} \wedge \text{state}_t = \text{passive} \}$
if $t > 0$
 then $t := t \ominus 1$
 else $\text{tour} := \text{second} ; t := n \ominus 1 ; tc := \text{white}$

T2: $\{ \text{tour} = \text{second} \}$
if $\text{color}_t = \text{black}$ **then** $tc := \text{black} ; \text{color}_t := \text{white} ;$
if $t > 0$
 then $t := t \ominus 1$
 else if $tc = \text{white}$
 then send $\langle \text{term} \rangle$ to all A_p
 else (* Reinitialize *)
 $\text{tour} := \text{first} ; t := n \ominus 1$

A_p: $\{ \text{state}_p = \text{active} \}$
 $\text{state}_q := \text{active} ; \text{color}_q := \text{black}$

I_p: $\{ \text{state}_p = \text{active} \}$
 $\text{state}_p := \text{passive}$

Algorithm 13: The ring-based termination detection algorithm.

and white here—this conforms to the usual description of similar termination detection algorithms. The color of the token (on its second tour) is stored in the variable tc .

The algorithm is initiated by process 0 by sending the token on its first tour to process $n \ominus 1$. A token visit during the first tour is described in action **T1**. It is enabled only when the process holding the token is *passive*, and consists of forwarding the token only (decrement t). If the token is at the end of the first tour it is whitened and sent on its second tour. A token visit during the second tour is described in action **T2**. If the color of the visited node is black, the token is colored black, and the color of the process is reset to white. While the token is not yet at the end of the second tour it is forwarded. At the end of the second tour, if the token is white termination is concluded, otherwise it is sent on its first tour again. Action **A_p** describes the (synchronous) activation of process q by process p and q 's subsequent blackening. Action **I_p** describes how process p becomes *passive*.

During the first tour, the token visits a process only in *passive* state. Thus a process can be *active* “behind” the token only if it is reactivated after the visit, but this implies that the process is black. Black processes are reported in the second round, regardless of what the basic computation does. The principle of

the algorithm is captured in the following predicate.

$$\begin{aligned}
P \equiv & \\
& \text{tour} = \text{first} \quad \wedge \quad (\forall q > t : \text{state}_q = \text{passive} \vee \exists q : \text{color}_q = \text{black}) \\
\vee & \text{tour} = \text{second} \quad \wedge \quad (\forall q : \text{state}_q = \text{passive} \vee \\
& \quad [tc = \text{black} \vee \exists q \leq t : \text{color}_q = \text{black}])
\end{aligned}$$

Lemma 3.7 *P is an invariant of the algorithm.*

Proof. After initialization $\text{tour} = \text{first}$ and $t = n \Leftrightarrow 1$ so P holds. It is easily verified that each of the actions maintains P. \square

Theorem 3.8 *Algorithm 13 is a correct termination detection algorithm.*

Proof. First suppose process 0 signals termination. This happens (see action T2) when the white token visits process 0 on its second tour, and $\text{color}_0 = \text{white}$. From invariant P, all processes are *passive*. Hence the safety property holds.

Now suppose the termination condition holds. No more blackening of processes occurs, so the next complete second tour whitens all processes, and after the subsequent second tour termination is signaled. This proves that the liveness property holds. \square

3.4.4 Discussion of the Algorithm

It is interesting to compare this algorithm with the similar algorithm by Dijkstra *et al.* [9]. In that algorithm a process is blackened upon *sending* rather than upon *receiving* an instantaneously transmitted message. Note that the subformula “ $\forall q > t : \text{state}_q = \text{passive}$ ” of P is falsified when process $q > t$ is activated by some $p \leq t$. The consequence of blackening upon sending rather than receiving is, that a black process is certainly found *ahead of* the token in this case (while in the new algorithm the black process may be found behind the token). But if this is the case, the two tours can be replaced by a single tour, and indeed the algorithm by Dijkstra *et al.* uses one tour only.

An alternative transformation of Ben-Ari’s algorithm uses blackening upon sending. Indeed, sender and receiver cooperate atomically in the \mathbf{A}_p action, which marks Z in Ben-Ari’s algorithm, and this virtual marking can be flagged in the sender as well as in the receiver. With this modification the \mathbf{A}_p action would be the following, and the reader may easily verify that this action also maintains invariant P above.

$$\begin{aligned}
\mathbf{A}_p: \quad & \{ \text{state}_p = \text{active} \} \\
& \text{state}_q := \text{active}; \text{color}_p := \text{black}
\end{aligned}$$

Unfortunately, invariant P is not strong enough to *exploit* the advantage of blackening upon sending, like in [9]. We conclude that the transformation of Ben-Ari’s garbage collection algorithm yields a termination detection algorithm

```

T: {  $state_t = passive$  }
    if  $color_t = black$  then  $tc := black ; color_t := white ;$ 
    if  $t > 0$ 
      then  $t := t \ominus 1$ 
      else if  $tc = white$ 
        then send (term) to all  $A_p$ 
        else (* Reinitialize *)
           $tc := white ; t := n \ominus 1$ 

```

Algorithm 14: The combined token visit.

which is very similar to Dijkstra’s algorithm, but less efficient because it needs two control tours rather than one. It is possible, however, to optimize Algorithm 13 by combining the second tour of one cycle with the first tour of the next cycle. The combined action **T** describing the token visit (see Algorithm 14) is almost identical to **T2**, but is “lazy” like **T1** in the sense that only *passive* processes are visited. When the combined action is used, the color of the initially active processes must be initialized to *black*.

4 Conclusions

In this paper we have presented a transformation of garbage collection schemes into termination detection algorithms. Termination detection and garbage collection belong to a family of problems concerning the detection of stable properties in distributed systems. This family also includes deadlock detection [5] and global virtual time approximation [32]. Several transformations of the same spirit as the one described here have been given between members of this family. It was argued by Natarajan [23] that termination is a special case of communication deadlock. It was shown by Tel [32] that termination detection is a special case of global virtual time approximation, and Mattern *et al.* [20] presented a derivation from termination detection algorithms to global virtual time approximation algorithms. Even though in some cases a solution to a “simple” problem is derived from a solution to a more “complicated” problem, we think transformations of this kind are useful for several reasons.

1. The research in one problem area may have “overlooked” a solution or optimization, which is possibly found by a transformation from a solution to another problem. We have derived not only known (e.g., Algorithm 9), but also new termination detection algorithms; see Algorithms 6, 10, and 13.
2. The design and verification of algorithms may be simplified.

3. The transformations show that the underlying difficulties (e.g., “behind the back messages”) of seemingly different problems are in fact similar if not identical. Hence similar solutions and techniques apply.
4. The transformations are of a theoretical interest because they enhance our understanding of the structure and intrinsic difficulties of a problem.

The transformation of garbage collection algorithms into termination detection algorithms raises some further questions, which will briefly be addressed in the remainder of this section.

4.1 Other Garbage Collection Algorithms

Virtually all garbage collection schemes can be transformed into sensible termination detection algorithms. Here we only sketch two more transformations, the reader is invited to complete the details and to apply the transformation to other garbage collection schemes (e.g., the well known algorithm by Dijkstra *et al.* [10]).

The “classical” garbage collection scheme consists in suspending the execution of the basic program when memory becomes short and run the garbage collector (of the mark-and-sweep type) while the program is stopped. Compared to on-the-fly garbage collection, synchronization and cooperation between the basic program and the collector is much simplified. The transformation of such a garbage collection algorithm yields a “freezing” termination detection algorithm where no reactivations are possible while the algorithm checks for the termination condition. In fact, one of the first published termination detection schemes (by Francez [13]) was a freezing algorithm.

Steele [30] describes a mark-and-sweep on-the-fly garbage collection algorithm. In this algorithm, when a reference from a marked to an unmarked object is installed, the marker process must visit the marked object again. In our transformation this principle means that when process p installs a reference to Z (i.e., becomes *active*), the termination detection algorithm must visit p again. This requirement can be realized in various ways. One way is to use a token visiting reactivated processes. Then the processes must keep specific information in order to record the identities of processes they reactivated (e.g., a vector with one component for each process).

4.2 A Different Transformation

The transformation of a garbage collection scheme into a termination detection algorithm as described in Section 2 proved to be very useful—a number of interesting termination detection schemes resulted from its application. However, there exists a different transformation principle which is in some respects “dual” to the principle we used up to now. In this section we sketch that principle, the details, however, are omitted.

Each process p is transformed into a non-root object A_p . In addition, a single (virtual) root object R exists. It is assumed that initially only a single object A_0 is *active*, all other objects are initially *passive*. Throughout the computation the following equivalence must be guaranteed:

$$A_p \text{ is } \textit{active} \Leftrightarrow R \text{ has an } A_p\text{-reference.}$$

This equivalence can easily be realized when a reference counting algorithm is used. The local reference counter RC_p is incremented when process A_p is activated, and decremented when it becomes *passive*. Thus, the actual references of R are only virtual references, and the object R need not be implemented at all.

In order to detect termination, we want to maintain the following property:

$$A_0 \text{ is garbage} \Rightarrow \text{the termination condition holds.}$$

Using this property a (distributed) garbage collection algorithm can detect the termination condition when collecting A_0 . The property is maintained by the following two rules:

1. Each basic message from object A_p to object A_q contains an A_p -reference.
2. When object A_q receives a message containing an A_p -reference it inserts that reference if it has no reference to any object, otherwise it immediately deletes the A_p -reference. If A_q is activated, RC_q is incremented (i.e., a virtual A_q -reference is installed at R).

The reader may easily verify that no cycles are formed and that if an object A_p is *active* then there exists a reference path from A_p (and consequently also from R) to A_0 . In order to guarantee the liveness property, the following rules should be observed in addition:

3. When object A_p becomes *passive*, RC_p is decremented. (That is, the virtual A_p -reference in R is deleted.)
4. When the reference counter RC_p of object A_p drops to zero and A_p has an A_q -reference, that reference is deleted ("recursive freeing" as part of the virtual collection of the garbage object A_p).

When Lermen and Maurer's garbage collection algorithm is used to detect whether A_0 is garbage (see Section 3.1), the two messages $\langle \mathbf{cop} A_p \rangle$ and $\langle \mathbf{ack} A_p \rangle$ can be merged into a single message. This is the case because A_p sends a reference pointing to itself. The only effect of the subsequent execution (in arbitrary order) of the receive actions \mathbf{RA}_p and \mathbf{RC}_p is to increment $aR_p(o)$ which is used as a guard for action \mathbf{DR}_p . Since A_p (virtually) sends $\langle \mathbf{inc} A_p, A_q \rangle$ to itself, action \mathbf{RI}_p is executed locally in A_p when sending a basic message.

The algorithm resulting after removal of the manipulation of references is simple. Whenever an object becomes *active* or sends a basic message it increments a local counter. The counter is decremented when the object becomes

passive or when a decrement control message is received. Decrement messages are either sent by rule 2 or by rule 4. The FIFO property is not required because the channel for which it was necessary in the original garbage collection algorithm no longer exists. Termination is detected when the counter of A_0 drops to 0.

The resulting algorithm is already known; it is the LRC termination detection scheme (Algorithm 12) derived in Section 3.3.2— A_0 plays the role of the environment e , and RC_p plays the role of the local reference counter LRC_p .

4.3 Reverse Transformation

It was already observed by Tel *et al.* [34] that detecting the termination of the marking phase emerges as a natural subproblem in mark-and-sweep garbage collectors. It was shown that the choice of a particular termination detection algorithm has a major influence on the resulting garbage collection algorithm. The transformation considered in this section is different, here it is indicated how a termination detection algorithm can be transformed into a reference counting garbage collection scheme.

The aim of a reference counting algorithm is to collect an object o when all o -references (in objects) have been deleted and no more o -references are in transit (in copy messages). A similarity to the termination detection problem is observed when activity of objects is defined suitably. An object is defined to be *o -active* if it holds an o -reference and *o -passive* otherwise, and a message is called an *o -activation message* if it carries an o -reference. Under these definitions, an *o -passive* object becomes *o -active* only upon receipt of an *o -activation message*, and only *o -active* objects send *o -activation messages*. The behavior of the computation is according to Algorithm 1, so that the *o -termination condition*, defined as

no process is *o -active* and no *o -activation messages* are in transit

is stable and can be detected by a termination detection algorithm. Furthermore

there are no o -references \Leftrightarrow the *o -termination condition* holds.

To arrive at a reference counting garbage collection algorithm, a termination detection algorithm is superimposed on the o -reference handling. When the *o -termination condition* is established, o is collected. For each object a separate instance of the termination detection algorithm is executed concurrently.

Not all termination detection algorithms can be reasonably used in this construction. Several considerations must be taken into account.

Centralized versus Distributed Control. It is not a drawback if an algorithm is chosen in which one process plays a special role, such as initiating the algorithm. The object o itself is a natural candidate to play this role. The central object could however become a bottleneck if its intervention is needed in every basic communication (as in the Activity Counting or Generational termination detection algorithm).

No Probe-Based Algorithms. The set of objects can be very large and varies due to creation and collection of objects. Therefore it is not feasible to use a termination detection algorithm in which *all* processes in the system take part. Rather, the activity of the algorithm should be restricted to processes that take part in the basic computation also. The LRC termination detection algorithm, the Credit Recovery, Activity Counting, and Generational termination detection algorithm all have this property, but probe-based algorithms are ruled out.

Early Termination. A process p holding or having held an o -reference may detect that it is garbage itself and must be collected. Therefore the termination detection algorithm must allow processes to terminate locally even while the computation as a whole has not yet terminated. The algorithm of [11] does not have this property: an “engaged” process must remain in the system as long as any of its descendants remain *active*.

These considerations differ from those that are usually taken into account when a termination detection algorithm is designed. For example, it is usually preferred that processes participate in the termination detection procedure only while they are *passive*, but this property probably conflicts with the possibility of early termination. The termination detection algorithms that we have derived from reference counting schemes present themselves as candidate algorithms, but their transformation yields no new algorithms, as the two transformations are each other’s inverse. Current research addresses the design of new termination detection algorithms that can more easily be used in the transformation sketched above.

4.4 Related Problems

The termination detection problem is an instance of a class of detection problems in distributed systems. Communication deadlock detection is a generalization where a part of the network can be terminated [23]; distributed infimum approximation is a generalization where the “property” to be detected takes values from any partially ordered domain, rather than just *passive* or *active* [32, 19, 20].

Deadlock Detection. In the communication deadlock problem, for each *passive* process a subset of the processes is determined at the moment it becomes *passive*. The process can become *active* only by receiving a message from a process in this subset. The termination detection problem is obtained, when each process always chooses the full set of processes. Thus, both the garbage collection problem and the communication deadlock detection problem seem to “dominate” the termination detection problem. This raises the question whether our approach can be generalized to detection of communication deadlocks.

\mathbf{S}_p : (* Send basic message *)
 send a message (M, x_p)

\mathbf{R}_p : { A basic message has arrived }
 receive message (M, x) ; $x_p := x_p \wedge x$

\mathbf{I}_p : (* Internal increase of x *)
 { $x_p < x$ } $x_p := x$

Algorithm 15: The basic actions of distributed infimum approximation.

Distributed Infimum Approximation. In the distributed infimum approximation problem an arbitrary partially ordered domain X with the *infimum* operator \wedge replaces the two-valued domain $\{active, passive\}$. The “state” x_p of process p is a value from X , and the messages of the basic computation are tagged with values from X . The handling of message tags and states in the operations of the basic computation are given in Algorithm 15. The problem is to approximate the global state function F , defined as the infimum of all states and tags of messages. It follows from the basic actions that this function is monotonically increasing.

Algorithms for this problem can be derived from termination detection algorithms, as described by Schoone and Tel [31] and Mattern *et al.* [20]. It would be interesting if our current construction could be generalized to obtain Distributed Infimum Approximation algorithms. To this end, instead of the single object Z a directed graph G_X of objects could be defined, reflecting the structure of X . The actions of the basic computation must then be formulated as reference manipulation, such that the growth of F is reflected by objects of G_X becoming garbage.

4.5 Legal Aspects

We have shown that termination detection algorithms are obtained as suitable instantiations of garbage collection schemes. Supplying a particular scheme, our transformation yields a particular termination detection algorithm. It may happen, however, that the resulting algorithm was found independently already. This is the case for example with the weighted reference counting scheme, which yields the Credit Recovery algorithm for termination detection, see Section 3.2.

Interestingly, commercial use of the weighted reference counting scheme by Watson and Watson is protected by a patent [12] which describes the invention in “a computer system having storage means containing memory cells, at least some of which contain pointers to others”. Does the patent now cover the Credit Recovery algorithm?

In recent years more papers described general transformations of solutions to one problem into solutions to another problem. The implication a patent in such a situation can have in general is too complicated for us and we are glad to leave it as food for lawyers [27].

Acknowledgements: We want to thank Martin Rudalics and Jörg Richter for their discussions of the paper and numerous suggestions and comments. We also thank the referees, whose comments have been a great help to improve the quality of this paper.

References

- [1] Ben-Ari, M., *Algorithms for On-the-fly Garbage Collection*, ACM Trans. on Prog. Lang. and Systems 6 (1984) 333–344.
- [2] Bevan, D.I., *An Efficient Reference Counting Solution to the Distributed Garbage Collection Problem*, Parallel Computing 9 (1989) 179–192.
- [3] Chandy, K.M., L. Lamport, *Distributed Snapshots: Determining Global States of Distributed Systems*, ACM Trans. on Computer Systems 3 (1985) 45–56.
- [4] Chandy, K.M., J. Misra, *How Processes Learn*, Distributed Computing 1 (1986) 40–52.
- [5] Chandy, K.M., J. Misra, L.M. Haas, *Distributed Deadlock Detection*, ACM Trans. on Computer Systems 1 (1983) 144–156.
- [6] Charron-Bost, B., G. Tel, F. Mattern, *Synchronous and Asynchronous Communication in Distributed Computations*, Technical Report, Université Paris 7, Paris, 1991.
- [7] Collins, G.E., *A Method for Overlapping and Erasure of Lists*, Comm. ACM 3 (1960) 655–657.
- [8] Chandrasekaran, S., S. Venkatesan, *A Message-Optimal Algorithm for Distributed Termination Detection*, Journal of Parallel and Distributed Computation 8 (1990) 245–252.
- [9] Dijkstra, E.W., W.H.J. Feijen, A.J.M. van Gasteren, *Derivation of a Termination Detection Algorithm for Distributed Computations*, Inf. Proc. Lett. 16 (1983) 217–219.
- [10] Dijkstra, E.W., L. Lamport, A.J. Martin, C.S. Scholten, E.F.M. Steffens, *On-the-fly Garbage Collection: An Exercise in Cooperation*, Comm. ACM 21 (1978) 966–975.
- [11] Dijkstra, E.W., C.S. Scholten, *Termination Detection for Diffusing Computations*, Inf. Proc. Lett. 11 (1980) 1–4.

- [12] European Patent Office, *Garbage Collection in a Computer System*, European Patent Application no 86309082.5.
- [13] Francez, N., *Distributed Termination*, ACM Trans. on Prog. Lang. and Systems 2 (1980) 42–55.
- [14] Goldberg, B., *Generational Reference Counting: A Reduced-Communication Distributed Storage Reclamation Scheme*, ACM SIGPLAN Notices 24 (July 1989) 313–321.
- [15] Ichisugi, Y., A. Yonezawa, *Distributed Garbage Collection Using Group Reference Counting*, Technical Report 90-014, Department of Information Science, University of Tokyo, 1990. (To be published by World Science Publishing Company: Proc. of the Annual Workshop on Software Science and Engineering, Kyoto, Japan, December 1989)
- [16] Lermen, C.-W., D. Maurer, *A Protocol for Distributed Reference Counting*, ACM Conference on Lisp and Functional Programming, Cambridge, 1986, pp. 343–354.
- [17] Mattern, F., *Algorithms for Distributed Termination Detection*, Distributed Computing 2 (1987) 161–175.
- [18] Mattern, F., *Global Quiescence Detection Based on Credit Distribution and Recovery*, Inf. Proc. Lett. 30 (1989) 195–200.
- [19] Mattern, F., *Efficient Distributed Snapshots and Global Virtual Time Algorithms for Non-FIFO Systems*, Tech. Rep. SFB124-24/90, Kaiserslautern University, 1990.
- [20] Mattern, F., H. Mehl, A.A. Schoone, and G. Tel, *Global Virtual Time Approximation with Distributed Termination Detection Algorithms*, Tech. Rep. RUU-CS-91-32, Dept. of Computer Science, University of Utrecht, 1991.
- [21] McCarthy, J., *Recursive Functions of Symbolic Expressions and Their Computation by Machine*, Comm. ACM 3 (1960) 184–195.
- [22] Misra, J., *Detecting Termination of Distributed Computations Using Markers*, Proc. of the 2nd ACM Symp. on Principles of Distributed Computing, Montreal, Quebec, 1983, pp. 290–294.
- [23] Natarajan, N., *A Distributed Scheme for Detecting Communication Deadlocks*, IEEE Trans. on Software Engineering SE-12 (1986) 531–537.
- [24] Piquer, J., *Indirect Reference Counting: A Distributed Garbage Collection Algorithm*, in: E.H.L. Aarts, J. van Leeuwen, M. Rem (eds.), Proceedings Parallel Architectures and Languages Europe, vol. I, Lecture Notes in Computer Science 505, Springer-Verlag, 1991, pp. 150–165.

- [25] Rudalics, M., *Multiprocessor List Memory Management*, Technical Report RISC-88-87.0, Research Institute for Symbolic Computation, J. Kepler University, Linz, 1988.
- [26] Rudalics, M., *Implementation of Distributed Reference Counts*, Technical Report (forthcoming), Research Institute for Symbolic Computation, J. Kepler University, Linz, 1990.
- [27] Samuelson, P., *Should Program Algorithms be Patented?*, Comm. ACM 33 (1990), 23–27.
- [28] Shavit, N., N. Francez, *A New Approach to Detection of Locally Indicative Stability*, in: L. Kott (ed.), *Proceedings ICALP 1986*, Lecture Notes in Computer Science 226, Springer-Verlag, 1986, pp. 334–358.
- [29] Van de Snepscheut, J.L.A., “*Algorithms for On-the-fly Garbage Collection*” *Revisited*, Inf. Proc. Lett. 24 (1987) 211–216.
- [30] Steele, G.L., *Multiprocessing Compactifying Garbage Collection*, Comm. ACM 18 (1975) 495–508.
- [31] Schoone, A.A., G. Tel, *Transformation of a Termination Detection Algorithm and its Assertional Correctness Proof*, Tech. Rep. RUU-CS-88-40, Dept. of Computer Science, University of Utrecht, 1988.
- [32] Tel, G., *Distributed Infimum Approximation*, Tech. Rep. RUU-CS-86-12, Dept. of Computer Science, University of Utrecht, 1986.
- [33] Tel, G., *Total Algorithms*, Algorithms Review 1 (1990) 13–42.
- [34] Tel, G., R.B. Tan, J. van Leeuwen, *The Derivation of Graph Marking Algorithms from Distributed Termination Detection Protocols*, Science of Computer Programming 10 (1988) 107–137.
- [35] Weng, K.-S., *An Abstract Implementation for a Generalized Dataflow Language*, Technical Report MIT/LCS/TR-228, Massachusetts Institute of Technology, 1979.
- [36] Watson, P., I. Watson, *An Efficient Garbage Collection Scheme for Parallel Computer Architectures*, in: J.W. de Bakker, A.J. Nijman, P.C. Treleaven (eds.), *Proceedings Parallel Architectures and Languages Europe*, vol. II, Lecture Notes in Computer Science 259, Springer-Verlag, 1987, pp. 432–443.

Contents

1	Introduction	2
1.1	The Termination Detection Problem	2
1.1.1	Description of the Problem	3
1.1.2	Solutions to the Problem	4
1.2	The Distributed Garbage Collection Problem	5
1.2.1	Description of the Problem	5
1.2.2	Solutions to the Problem	7
2	Termination Detection Using Garbage Collection	8
3	Examples of the Transformation	10
3.1	Distributed Reference Counting	10
3.1.1	Description of the Scheme	11
3.1.2	Transformation into a Termination Detection Algorithm .	13
3.1.3	Discussion of the Algorithm	14
3.1.4	Related Algorithms	16
3.2	Weighted Reference Counting	18
3.2.1	Description of the Scheme	18
3.2.2	Transformation into a Termination Detection Algorithm .	19
3.2.3	Discussion of the Algorithm	21
3.3	Local Reference Counting	23
3.3.1	Description of the Scheme	23
3.3.2	Transformation into a Termination Detection Algorithm .	24
3.3.3	Discussion of the Algorithm	26
3.4	Mark-and-Sweep Garbage Collection	26
3.4.1	Description of the Scheme	27
3.4.2	Transformation into a Termination Detection Algorithm .	27
3.4.3	Correctness of the New Algorithm	28
3.4.4	Discussion of the Algorithm	30
4	Conclusions	31
4.1	Other Garbage Collection Algorithms	32
4.2	A Different Transformation	32
4.3	Reverse Transformation	34
4.4	Related Problems	35
4.5	Legal Aspects	36