

Recent Advances in Distributed Garbage Collection

Marc Shapiro^{1,3}, Fabrice Le Fessant^{1*}, and Paulo Ferreira²

¹ INRIA Projet SOR, Rocquencourt, France,
<http://www-sor.inria.fr/>

² INESC, Lisboa, Portugal

³ Microsoft Research Ltd., Cambridge, United Kingdom

1 Why Distributed Garbage Collection

Dynamically-allocated memory must eventually be reclaimed. But manual reclamation is error-prone, and if an object is de-allocated prematurely, a program that later follows a pointer to it might behave incorrectly. In contrast, *garbage collection* (GC) automatically reclaims objects that can no longer be reached by any path of pointers.

In a distributed system, an object might not only be referenced from one program, but also from other programs and other computers. This makes manual reclamation quite intractable. We consider here the problem of *distributed garbage collection* (DGC). DGC is particularly important in a large-scale distributed system, since a de-allocation error might cause a completely unrelated program to fail, possibly far away and an arbitrarily long time later.

DGC has been a subject of academic research since at least 1977 [4]. With the growth of the Internet, DGC is now receiving its share of commercial attention. For instance, both Java RMI [28] and DCOM [23] come with some form of DGC.

The next section contains a quick review of DGC techniques; for a more in-depth treatment, we refer to published surveys [1, 21]. Then we detail two recent advances in DGC, both developed in the context of the Broadcast project. One is an algorithm for collecting distributed cycles of garbage in a message-passing system; the other is an algorithm for DGC in a system with caching and/or replication. Both have been implemented in real systems and are in actual use.

This article assumes some familiarity with centralised GC, now a mature area [11, 27]. We use the following vocabulary. Objects connected by references form a graph; a remote reference may point into another process. An application or mutator enters the graph via roots, allocates objects, and performs assignment of reference variables. The system or collector reclaims garbage, i.e., objects that are reachable by no path from any root.

2 Distributed Garbage Collection Approaches

Distributed algorithms should be scalable, which implies local execution, absence of remote synchronisation, low complexity and costs, and fault tolerance.

* Ph.D. student at École Polytechnique. Also with INRIA projet PARA.

Scalability in DGC is often achieved at the cost of incompleteness, i.e., only reclaiming a safe subset of actual garbage. Scalability excludes the apparently straightforward approach of running one of the existing centralised algorithms over a consistent view of the object graph, e.g., over a snapshot [5] or in a transaction. Our approach instead is to use the unique properties of DGC that enable specific, efficient solutions.

2.1 Two Models of Distributed Systems

Before going into more detail, it is appropriate to explain our two models of distributed system. Classically, a distributed system is defined as a set of disjoint processes, sharing no memory, and communicating by *message passing*; this is our first model. We do not assume any upper bound on message delivery time. An object exists at a single process at any time.

In the alternative *shared-memory model*, processes interact by sharing memory. The system simulates a shared memory by sending (in a message) the value of a chunk of memory to a remote site, where it is cached and mapped into application space. We call the different copies of a given chunk of memory its replicas. The system keeps replicas consistent, detecting updates and sending messages invalidating or updating out-of-date replicas. Messages are not visible to applications. This model includes Distributed Shared Memories (DSMs) such as Ivy [18], distributed shared object systems such as Orca [3], and cached distributed file systems or databases.

We use the word *bunch* for the smallest subdivision of the shared memory that can be replicated, often a page or segment. An object resides in a single bunch, but at any point in time the bunch (and the objects it contains) might be replicated to multiple processes.

2.2 Basic DGC Techniques

The simplest DGC algorithm uses time-outs, as was first proposed in the Bullet file system [26]. Time-outs are simple and effective but unsafe, since a transient error can cause a reachable object to disappear.

Other DGC algorithms derive from the two well-known families of GC algorithms, namely *counting* and *tracing*. A counting algorithm maintains a count of references to every object. Counting is a local operation, hence the algorithm scales well. It has the disadvantage that it does not collect cycles of garbage. *Reference listing* generalises counting to maintain the list (and not just the number) of pointers to an object; this allows a reference originating from a site declared crashed to be reclaimed.

A tracing algorithm determines garbage by direct examination of the graph. A tracing collector collects all garbage, including cycles. However distributed tracing is not scalable, because it is a global algorithm, and because it runs globally synchronised phases.

Some authors advocate *back-tracing* [20]: if a backwards graph walk returns to the start object and does not encounter a root, then that object is part of a

garbage cycle. The backwards walk is expensive; therefore we do not consider back-tracing any further.

2.3 Hybrid Collection Algorithms

Practical DGC algorithms are a hybrid of tracing, counting, and time-out. For instance a hybrid DGC is used in the commercial remote-object system Java RMI [28]. Two systems that we developed, SSP Chains [24] and PerDiS [9], are also based on hybrid DGC; they will be examined in more detail in the later sections.

A distributed system is typically partitioned for locality. In the DGC context, we call each part a *space*. DGC algorithms typically focus on collecting references that cross the boundaries between spaces. Such remote references are managed with reference lists.

Each space independently takes care of its local garbage with a tracing algorithm.¹ An incoming remote reference is (conservatively) considered part of the local root set. When a local collector determines that an outgoing remote reference has become unreachable, the DGC sends a message to the target space, causing it to decrement its reference count. Time-outs are used to detect that a space has crashed; reference listing allows to selectively ignore references from a crashed space.

A hybrid DGC combines the best of all worlds. Within a space, the local collector will reclaim all garbage, including cycles, without imposing a global algorithm. Across spaces, listing provides scalability.

However, a hybrid algorithm does not collect cycles of garbage that span spaces. Furthermore, time-outs must be used with utmost care.

Hybrid Collection in the Message-Passing Model. In the message-passing model, each process (or processor) is a space. It is easy to detect reference mutations, because assigning a remote pointer requires a message, containing the reference, between processes. Then, sending or receiving the message causes the incrementation of the corresponding count.

This algorithm model is well understood; remaining open questions have to do with fault tolerance and with collecting distributed cycles of garbage. Section 3 presents a new algorithm for collecting distributed cycles in a large-scale distributed system.

Hybrid Collection in the Shared-Memory Model. In the shared-memory model, each bunch replica constitutes a separate space. At some arbitrary point in time, replicas of a same bunch may be mutually inconsistent. Multiple spaces may coexist in a single process. A pointer to an object designates any of its replicas, independent of process. Pointer assignments occur directly in memory.

¹ Local tracing is not an obligation; local counting or even manual collection is possible too.

The collector must somehow be notified when the application assigns a pointer across a bunch boundary. This all makes it hard to determine which objects are reachable. The advance described in Section 4 is the Larchant algorithm, an efficient, scalable DGC algorithm for a distributed shared memory, which takes these issues into account.

3 Collecting Distributed Garbage Cycles in the Message-Passing Model

In this section, we present our recent work on collecting distributed cycles of garbage in message-passing systems. This work is mainly inspired from Hughes' algorithm, but our detector [15] extends Hughes to asynchronous distributed systems.

After introducing the main properties of our detector, we will present the algorithm and argue for its scalability and fault-tolerance. More detail can be found elsewhere [15].

3.1 Introduction

Terminology. In this paper, we use the terminology of the Stub-Scion Pair Chains system (SSPC) [24].

Stubs and Scions. Each remote reference R from object A in space X to object B in space Y is represented by a local pointer in X from object A to a special object $stub_X(R)$, called a *stub* and used as a proxy in X for object B , and a local pointer in Y from another special object $scion_Y(R)$, called a *scion* and used as a local root in Y , to the object B . For the reference R , X is the *upstream* space, and Y the *downstream* space.

Stubs and scions come in pairs: each stub has exactly one *matching* scion, and each scion has at most one matching stub. A remote reference R is created by first creating the scion, then sending the reference (in fact, the scion identifier, called the *locator*) in a message, and finally by creating the associated stub in the remote space.

Timestamps and Dates. In the SSPC acyclic garbage collector, all messages are stamped to prevent race conditions. In our algorithm, we use both the original SSPC timestamps, and other timestamps (starting times of global traces, generated by a distributed Lamport clock) on stubs and scions to distinguish multiple concurrent traces. To avoid confusion, the former are called *timestamps*, while the latter ones are called *dates*.

Overview. Like Hughes' algorithm [10], our detector is based on multiple global traces progressing concurrently in the system. Each global trace starts as some space's local garbage collection. It uses as its marker the current date from the

Lamport clock at its starting space, and propagates that date from local roots to stubs reachable from those roots. Dates on stubs are then propagated along chains of remote pointers, from stubs to scions by messages and from scions to reachable stubs by local garbage collections. A stub is always marked with the highest date of all scions or local roots it is known to be reachable from.

As a consequence, reachable stubs are marked with increasing dates, remotely-propagated from their original roots. On the contrary, unreachable stubs eventually stabilise at the starting date of the latest global trace when they were found to be reachable.

A global trace is *terminated* when its starting date no longer appears on reachable stubs and scions. Using an analysis of increasing dates on its stubs, a space computes the date (noted *localmin*) of the oldest global trace which has not yet terminated at this space. All *localmins* are gathered on a central location, called the *detection server*, in order to compute the oldest date (noted *globalmin*) among traces which have not terminated in one or more spaces. Any scion that is marked with a date earlier than *globalmin* is removed from the roots, since it is marked with the date of a terminated trace.

Main Properties. This algorithm presents interesting properties for distributed systems:

- *It is centralised:* centralisation is often seen as a drawback in distributed systems. However, for this problem, it enables a simpler, lightweight and fault-tolerant solution. We will show that the central server is not a performance or fault-tolerance bottleneck.
- *It is topology-aware:* A network is often organised as a closely-coupled LAN, loosely connected to the rest of the world through a WAN. In our model, there would one server per LAN, and only intra-LAN cycles are detected. Inter-LAN cycles are not detected.
- *It is optional:* SSPCs constitute the basic reference and communication mechanism. The new algorithm presented here detects detects cycles of garbage spanning participating spaces. A process may decline to participate while still using SSPCs normally. A cycle of garbage that goes through a non-participating space is not collected, but there is no further overhead for the non-participant.
- *It is asynchronous:* messages are sent asynchronously. Neither a mutator nor a collector is ever blocked by the DGC.
- *It consumes few resources:* our algorithm relies on small local data structures, it sends few messages, only to spaces in the immediate vicinity, and its computation overhead is negligible.
- *It has a low implementation cost:* our algorithm has been implemented for the Objective-CAML [17] platform of Stub-Scion Pair Chains [13]. Only minor modifications of the runtime were necessary to propagate timestamps as a side-effect of local garbage collection. Other parts of the algorithm were implemented over the runtime as an optional library.

- *It is fault-tolerant*: the algorithm is safe in the presence of unreliable communications (messages can be lost, re-ordered or duplicated) and crashes.

Comparison with Hughes’ Algorithm. Hughes’ assumptions are quite stronger than ours: his algorithm relies on a global clock, on instantaneous and reliable communication, and does not tolerate space crashes. Thus, important aspects of large-scale distributed systems, such as messages in transit or failures, are not addressed. These strong assumptions enable to detect termination using Rana’s algorithm [22], based on a snapshot of local states.

Our algorithm replaces Rana’s algorithm by the centralised computation of the minimum of an integer vector. Our assumptions are more realistic: communications are assumed asynchronous and unreliable, and the detection process is not broken by space crashes. Hughes’ global clock is replaced by a Lamport logical clock [12]. Moreover, most computation is delayed until local garbage collection time, leading to a more conservative but cheaper algorithm.

3.2 The Algorithm: Propagation of Dates

Local Propagation. At the beginning of each local garbage collection, local roots are marked with the current date of the local Lamport clock. Local roots and scions are then sorted,² and traced in decreasing order of their dates. A reachable stub is marked with the date of the root being traced. Since an object is only marked once during a local collection, a stub ends up marked with the greatest date of all roots from which it is reachable.

Remote Propagation. After a local collection, a marked stub whose date has increased, propagates its date to its downstream scion, using a **Stubdates** message. A stub whose date has remained constant or has decreased does not propagate anything.

SSPC timestamps are used to check if a scion locator has been sent, and not received before the **Stubdates** message was sent, in which case the scion date³ must not be updated, since the new reference was not taken into account in the propagated dates.

Characterising a Cycle of Garbage. At each local garbage collection, local roots propagate a new larger date. Thus, the date of any reachable stub will eventually increase. In contrast, dates belonging to unreachable cycles evolve in two phases: the greatest date is first propagated to all stubs and scions in the cycle. Then, since there is no local root leading to the cycle, no new date can enter the cycle, and the dates remain constant forever (see figure 1).

² Scions are marked with the date received from their associated stub. Initially, they are marked with a special date, called *top*, which is always replaced by the current date in computations.

³ The scion date is always set to *top* when its locator is sent.

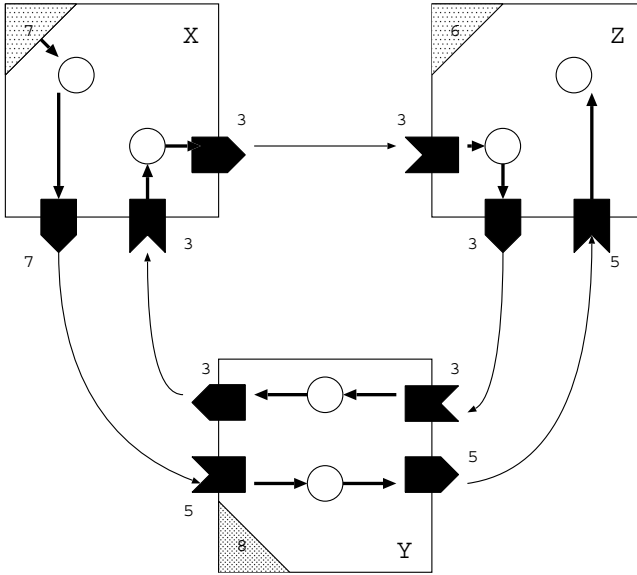


Fig. 1. As a consequence of date propagation, unreachable stubs in a cycle are eventually marked with the date of a terminated trace (date 3), whereas dates on reachable stubs continue to increase with successive global traces (dates 5 and 7).

The system is recurrently traversed by concurrent traces originating at the different local roots. Each one is characterised by its marker, the date at which it started. A reachable stub is recurrently marked by a new trace, i.e., with a greater marker. An unreachable stub, in contrast, remains marked with the the starting date of the last trace that found it to be reachable.

The goal of our algorithm is therefore to compute an increasing threshold, called **globalmin**, which is a minimum on dates found on stubs marked by successive global traces. **Globalmin** is a conservative approximation of the date of the oldest concurrent trace that has not yet terminated somewhere.

3.3 The Algorithm: Computation of **globalmin**

Computation of **globalmin by the Detection Server.** The computation of **globalmin** requires a consensus involving all participating spaces. A fully distributed and fault-tolerant consensus would be extremely complex. However, for our particular case, the use of a centralised computation will be shown to be scalable and fault-tolerant.

In our system, **globalmin** is computed by a dedicated process, called the *detection server*. Each participating space computes a value (called **localmin**), rep-

resenting the date of oldest trace that has not yet terminated at this space. The space then sends **localmin** to the detection server in a **Localmin** message. Since a trace that is locally terminated can be revived by receiving the marker of that date from another space, the value of **localmin** varies non-monotonically.

However, **globalmin**, the conservative minimum of all **localmins**, is monotonically increasing, since a trace is terminated only when terminated for all spaces. An unreachable stub or scion is marked with the date of a terminated trace, which eventually becomes smaller than **globalmin**.

The value of **globalmin** on the server does not need to be always up-to-date. Keeping an old value of **globalmin** means that the termination of some traces has not been detected by the detection server yet, and therefore that the reclamation of cycles with those dates is simply delayed. To reduce the load on the detection server, **globalmin** can be computed unfrequently, with a period depending only on the actual need of memory space in the participating spaces.

The only information maintained by the detection server is the list of participating spaces and, for each participating space, the last value of **localmin** received. This information is simple to recover in the case of a crash of the detection server. A new detection server can be started promptly in the same computation state, without aborting any global trace.

Computation of localmin by a Participating Space. For any participating space, **localmin** is a conservative approximation of the oldest date of the traces that have not terminated yet for that space. We say that a space *protects* a date, preventing any stub or scion marked with that date from being reclaimed, when its computed **localmin** is smaller than, or equal to, that date.

The value **localmin** is computed according to the following *protection* rule:

- When the local garbage collector of some space increases the date of a stub, that stub’s previous date must remain protected by the space until the collector of the downstream space has itself propagated the new date from the stub’s downstream scion, and the resulting **localmin** (computed by the downstream space) has been received by the detection server (see Figure 2).

Indeed, if the local garbage collector increases the date of a stub, that stub is probably reachable. Since the stub’s downstream scion is still marked with the previous date of the stub, the global trace for that date is not terminated yet.

However, the downstream space will not protect that date until it has itself detected that the associated global trace has not terminated yet. This detection will occur when the increased date has been propagated to the downstream space by a message, and to possibly reachable stubs by a local garbage collection. Finally, for the protection to be effective, the detection server must have received the associated **localmin** message from the downstream space.

This rule can be implemented by the following protocol: The local garbage collector computes, for each downstream space, the minimum date which must be protected, called **ProtectNow**, by comparing the new and old dates on stubs. **ProtectNow** entries are entered into the **ProtectedSet** of all dates protected on

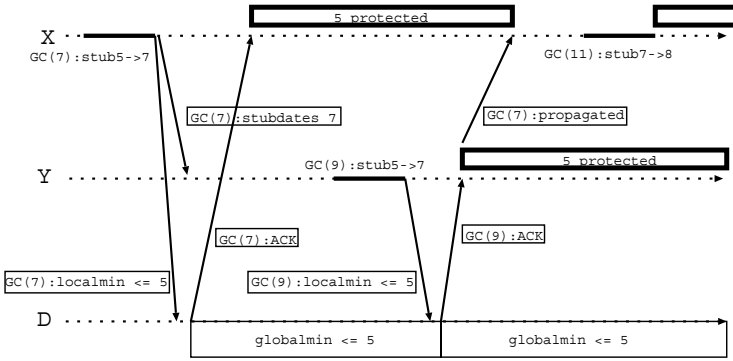


Fig. 2. Our protection rule is safe: space X stops protecting date 5 when Y has started protecting it.

the space. **Localmin** is then computed as the minimum of all entries in the **ProtectedSet**.

Localmin is sent to the detection server, which returns an **Ack** message, containing the current value of **globalmin**. When the **Ack** is received, the space knows that all new dates received before the local garbage collection have been propagated and that the new **localmin** has been received by the server. According to the protection rule, upstream spaces may stop protecting old dates associated with the stubs and local garbage collection whose new dates have just been propagated.

Moreover, when a space receives a **Stubdates** message, it stores the associated upstream local garbage collection date in a set, called the **PropagatedSet**. When an **Ack** message is received, the space sends a **Propagated** message to each upstream space mentioned in **PropagatedSet**. This message contains the entries of the **PropagatedSet** entered before the acknowledged local garbage collection. When an upstream space receives **Propagated**, it may remove the entries in its **ProtectedSet** for the sending space for all local collections up to the propagated one.

Coping with Mutator Activity. The mutator may send remote references (thus creating stub-scion pairs or using existing ones), or invoke reference targets, while the distributed collector runs.

When a new stub-scion pair is created, the scion date is initialized to *top* and the stub date to the current date of its space. When re-using an existing stub-scion pair, the scion owner does not know if the stub associated with its scion is still live. Thus, it must behave as if a new pair were created: the scion date is set to *top*. The *top* date on a scion is only replaced when a new date is propagated from its matching stub. However, to avoid race conditions, the

Stubdates message propagating the new date must have been sent after the last message containing the scion locator has been received by the sender space. Observe that this simple mechanism also handles non-participating spaces, since scions reachable from them remain marked with *top*, and thus are never collected.

When the mutator invokes an object through a stub-scion pair, this is sufficient proof that the stub is reachable. Therefore, an invocation increases the stub's date to the current date. As a consequence of the protection rule, its previous date must therefore be protected by subsequent **localmin** values.

3.4 Complexity

We can now examine the complexity of our algorithm, in memory, computation time and communication requirements.

Participating Spaces

Memory Consumption. Compared to the basic SSPC, each stub contains two extra dates and each scion one. Our algorithm needs two new data structures at each participating space. **ProtectedSet** contains the dates to be protected for each remote space by the value of **localmin**. **PropagatedSet** contains the dates of **Stubdates** messages received before a local garbage collection for which **Localmin** message is awaiting acknowledgement. Both can be implemented as FIFO queues, one for each remote participating space.

One entry is put in each queue of **ProtectedSet** at each local garbage collection. Therefore size of a queue depends on the frequency of local garbage collections with respect to remote ones. If the frequencies are similar across spaces, each queue should contain a small number of entries. However, if a space collects more frequently, causing its set to become too large, date propagation and computation of **localmin** should be avoided until sufficient entries have been removed from its **ProtectedSet**.⁴

The size of the **PropagatedSet** depends on the delay between transmitting a **Localmin** message and receiving the **Ack**. Thus, it should be small unless communication is very unreliable.

Computation Time. The major computation load on a space occurs during local garbage collection. The local trace requires that scions are sorted by decreasing order of their dates. This sort can be done from scratch (cost $O(N \log(N))$ in the number of scions) at each local garbage collection, or by inserting a scion into a sorted table each time its scion date is modified. If there are fewer dates than scions, another approach is to group scions by date, and sort the groups.

Localmin is computed as the minimum of the old dates of increasing stubs, and of all the entries in the **ProtectedSet**. Thus, there is one comparison for

⁴ Notice that one remote garbage collection can remove several entries at once, i.e., all the entries for local garbage collections whose **Stubdates** messages were received before the date of the remote garbage collection.

each stub, to compute the current date to be protected for the current local garbage collection for each remote space, and one comparison for each entry in the `ProtectedSet`, i.e., $O(n_{stubs}) + O(n_{spaces})$.

Detection Server. The detection server, although centralised, is not a bottleneck. Indeed, the computation and memory requirements are low, and linear in the number of participating spaces:

Memory consumption: The detection server maintains a vector containing the last `localmin` received from each participating space.

Computation time: `Globalmin` is computed periodically as the minimum of the `localmins` of the participating spaces. To trigger this computation, a good stimulus is that all participating spaces whose previous `localmins` were equal to `globalmin` have increased their `localmin`. The number of such spaces can easily be computed during the computation of `globalmin`.

Messages. For each local garbage collection, two messages are sent to each remote space in the vicinity: `Stubdates` and `Propagated`. There are also two messages exchanged with the detection server: `Localmin` and `Ack`. Whereas `Propagated`, `Localmin` and `Ack` are small messages (containing at most three values), `Stubdates` contains one date per live stub. However, this message can be merged with the `Live` message of basic SSPC, which also contains one date per live stub.

3.5 Fault-Tolerance

This algorithm is tolerant both to unreliable communication and to space failures. Although building reliable communication above unreliable communication is often not a problem, space failures cannot be avoided in a real distributed system.

Unreliable communication is tolerated because of the conservative approach taken in the communication of `globalmin` and `localmin`. `Localmins` are computed not only from the protected dates of the current garbage collection, but also for some previous garbage collections (`Protected Set`). `Globalmin` is computed from the last `localmins` received, and its value is always safe, even if not up-to-date.

The detection server can crash. Any other processor can then become the new detection server, without aborting either the current detection process nor the current global traces.

A participating space can also crash. In such a case, the detection server must first propagate the new list of participants to all remaining ones. If the upstream space of a scion is not known to be one of the participants, its date is set to the special value *top*. Then, the detection server waits for the new `localmin` value computed after this update, before computing the next value of `globalmin`.

3.6 Scalability

Our algorithm imposes few requirements on participating spaces. The computational and memory overhead is low, and few messages are sent for each local garbage collection. However, it has two major drawbacks: the computation is centralised, and a slow space will slow down the detection of distributed cycles of garbage (but not the detection of non-cyclic garbage, which is down by the base SSPC system).

However, we consider this algorithm interesting in large scale networks. Indeed, the limitation introduced by these two drawbacks is on the *number* of participating spaces, and not the *distance* between them. Since cycle detection is most useful between long-lived (and/or persistent) servers, these limitations are not problematic. A short-lived client does not need to participate in cyclic collection, since any garbage cycles through it will be collected when it dies. Only long-lived processes should participate.

3.7 Conclusion

We have described a detector of distributed cycles of garbage. Our cycle detector presents some interesting properties for large-scale systems: asynchrony between participating spaces, optional participation, tolerance to communication faults and space crashes, low resource requirements, and ease of implementation (no modifications to local objects, only a few to the local garbage collector). Moreover, this algorithm has already been implemented in a distributed system, and the implementation details can be found in le Fessant [15].

We are now working on a new version of this cycle detector using both propagation of marks and back-tracing. This new algorithm will have the same properties as the one described here, plus the ability to detect cycles spanning any spaces in the whole system. This new algorithm is currently being implemented for our mobile agents platform [14].

4 The Larchant DGC Algorithm for Distributed Shared Stores

4.1 Introduction

Modern distributed systems rely on caching and replication to speed up access to remote objects. A reference can then be resolved directly in the local memory. The purest example is provided by a Distributed Shared Memory (DSM) [18], which enables processes on different computers to share data simply by using pointers, just like in a centralised program. Distributed file systems and object-oriented databases provide similar facilities. In what follows we will refer to all such cached or replicated memories as DSMs.

DGC in a distributed shared memory is a harder problem than DGC in the message-passing model studied in Section 3 because:

- Applications modify the graph concurrently and by pointer assignment, not by sending messages. This is a very frequent operation, which should not be slowed down; for instance reference counting at each pointer assignment would not be acceptable.
- Replicas are not instantly coherent. Observing a consistent image of the graph is difficult and costly.
- The pointer graph may be very large and distributed. Much of it resides on disk. Tracing the whole graph in one go is unfeasible.
- A single assignment can affect vast and remote portions of the graph. This has consequences on the global ordering of operations.
- GC should not compete with applications. For instance, it should not take locks, cause coherence operations, nor cause I/O.

We consider previous results on collecting a DSM [2, 16, 19, 29] inapplicable because it does not take the above issues into account. We describe now work done in the context of Larchant, a distributed and persistent shared store, intended for interactive cooperative tasks. Larchant consists essentially of a large-scale DSM; only objects reachable from a *persistent root* (e.g., a name server) persist.

The main goals of our distributed GC algorithm are correctness, scalability, low overhead, and independence from any particular coherence algorithm. Secondary goals are avoiding source code and compiler changes.

Our approach divides the global GC into small, local, independent pieces, that run asynchronously, hence can be deferred and run in the background:

- The store is partitioned into “bunches;” a bunch may be replicated. GC is a hybrid of tracing within a bunch and counting across bunches.
- Each site runs a collector with a standard tracing algorithm [27] that works in one or more bunch replicas (on that site) at the same time.
- The cooperation protocol between collectors does not entail mutual synchronisation.
- A collector examines only the local portion of the graph, without causing any I/O or taking locks.
- A collector may run even when local replicas are not known to be coherent.

We exhibit five simple rules for the correctness of collection in a DSM, which Larchant satisfies. The algorithm is safe (no reachable data is reclaimed) and live (garbage is eventually reclaimed). Sadly, but unavoidably, it is not complete.

The underlying theoretical model and a proof of correctness can be found in Ferreira and Shapiro [8]. More information about the Larchant algorithm can be found in previous publications [6, 7, 25].

4.2 GC Algorithm and Safety Rules

The Larchant memory, as shown in Figure 3, is subdivided into coarse-grain *bunches* (typically, a set of contiguous pages) in which objects are allocated. A bunch can be replicated (e.g., cached) in multiple processes. An object can

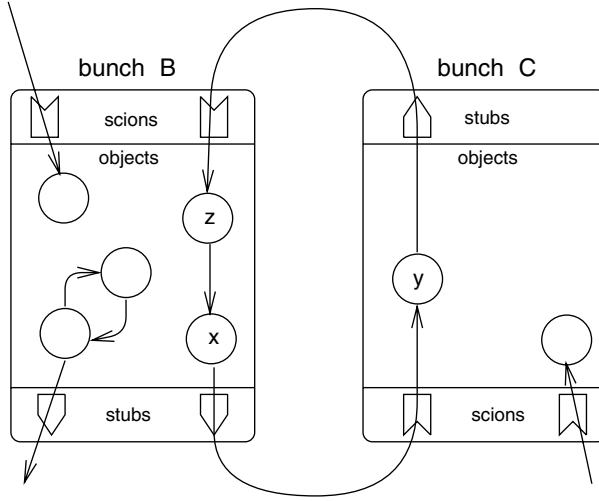


Fig. 3. Two bunches containing objects, stubs, and scions.

point to another object, either in the same bunch or across bunches. In both cases the reference is a raw pointer. In the second case only, the collector maintains additional data structures, a *stub* for an outgoing pointer and a *scion* for an incoming one.⁵ A typical execution example is shown in Figure 4.

Larchant uses a hybrid DGC, as defined in Section 2.3. Each replica of each bunch forms a space.⁶ A reference within a bunch is subject to tracing collection, whereas references that cross bunch boundaries are counted. To avoid interfering with the application, we assume that the collector discovers pointer assignments after the fact, as it traces a bunch. Tracing a bunch causes stubs to be created or deleted; the counting algorithm adjusts scions accordingly.

We now outline the basic algorithm and define some notation. When a mutator performs the assignment noted $\langle x := y \rangle_i$ (i.e., copy the value of pointer y into pointer x in process i), up to three processes are involved in the corresponding counting. Say objects x , y , z and t are located in bunches X , Y , Z and T respectively; prior to the assignment, x pointed to z and y pointed to t . As a consequence of the assignment, the collector of process i increments the reference count for t by performing the local operation $\langle \text{increment.stub}(Xx, Tt) \rangle_i$, and sending message $\text{increment.scion}(Xx, Tt)$ to the collector in some process j managing T . It also decrements the reference count for z by performing the local

⁵ Despite the similar vocabulary and role, these stubs and scions are somewhat different from those in SSPC used in Section 3. Larchant's stubs and scions are auxiliary data structures of the distributed collector; the mutator does not see them in any way; increments and decrements occur in the background.

⁶ To the first approximation only, because Larchant changes space boundaries dynamically by grouping together a number of bunches [25]. However this is beyond the scope of the current discussion.

operation `decrement.stub(Xx, Zz)`, and sending message `decrement.scion(Xx, Zz)` to the collector in some process k , managing Z . These adjustments of reference counts need not occur immediately, as we will see later.

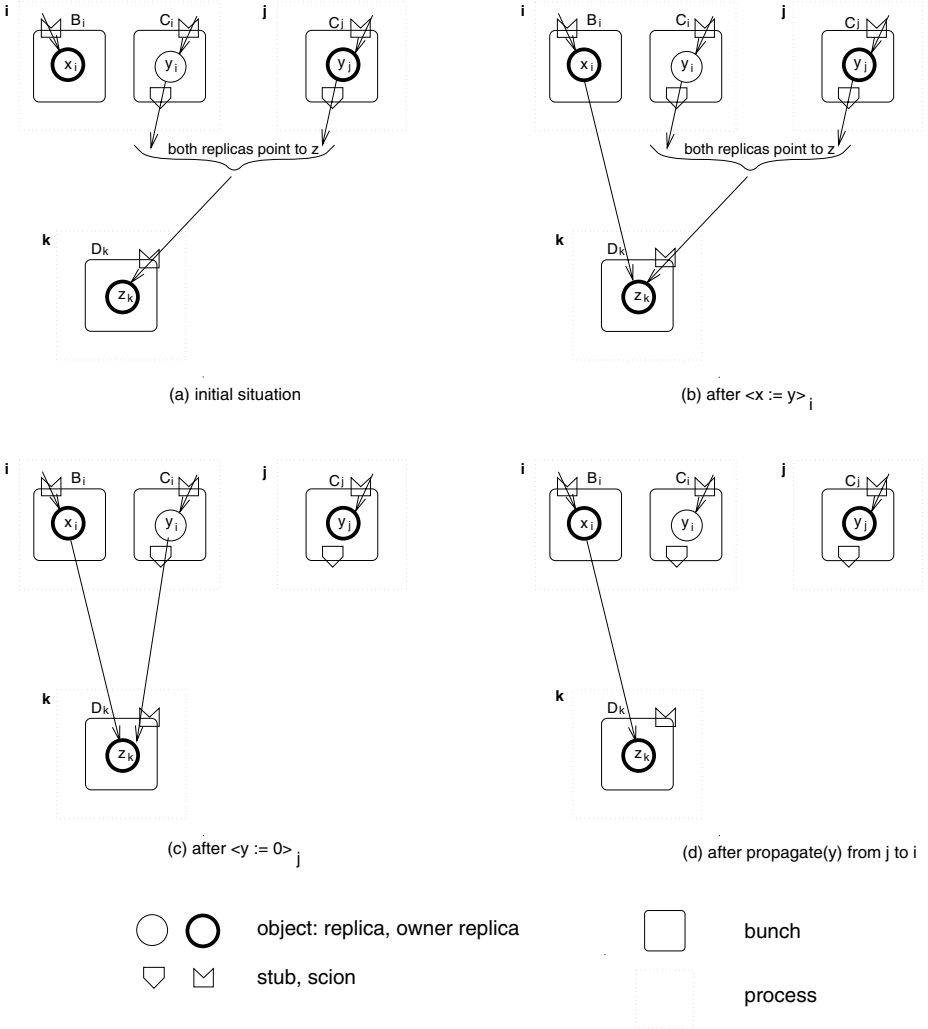


Fig. 4. Prototypical example of mutator execution. Note that the stubs and scions become temporarily inconsistent with the pointers. However, as described in the paper, this does not compromise safety.

Tracing in the Presence of Replicas: The Union Rule. Each process runs a standard centralised tracing collector. The issue we raise now is how collectors cooperate, in order to take replication into account. It is desirable that a collector remain independent, both of remote collectors, and of the coherence algorithm. Thus, a collector may scan a local replica, even if it is not known to be coherent, and independently of the actions of remote collectors.

The collector at process i might observe x_i to be pointing to z , whereas collector at process j concurrently observes x_j to be pointing to t . The coherence protocol will eventually make both replicas equal, but the collector cannot tell which value of x is correct. In the absence of better information, the collector must accept all replicas as equally valid, and never reclaim an object until it is observed unreachable in the union of all replicas. This is captured by the following rule.

Safety Condition I: Union Rule. *If some replica x_i points to z , and some replica x_j is reachable, then z is reachable.*

The above says that if some object z is referenced only by an unreachable replica x_i , it is reachable nonetheless if some other replica x_j , $i \neq j$, of x is reachable. This very conservative formulation is necessary in the absence of knowledge of the coherence algorithm.

An efficient implementation of the Union Rule applies to single-owner coherence protocols such as entry consistency. In such protocols, a single “owner” process centralizes the information about a given object.⁷ The collectors centralise the information about pointers from x at the owner of x , using what we call union messages. In other words, a process holding replica x_j sends a union message, to x ’s owner, after detecting a change in the pointers from x_j . (Note that this detection is achieved by tracing x_j ’s enclosing bunch.)

Now, suppose that x points to z , and x is assigned a new value (for instance the null pointer). It is only when all the replicas of x have the new value, and the corresponding collectors have informed x ’s owner (by sending it a union message) that there are no pointers from x to z , that the owner of x sends a message to the owner of z , to decrement the corresponding scion’s reference count. This technique moves some of the responsibility for reference counting to the owner of the objects where references originate.

Cross-Bunch Counting and More Safety Rules. The standard approach to reference counting is to instrument assignments in order to immediately increment/decrement the corresponding counts. This approach requires compiler modification, and is expensive when assignments are frequent and counting is a remote operation, as is the case in Larchant.

Our solution consists of deferring the counting to a later tracing. In fact, the counts need not be adjusted immediately. Consider an assignment $\langle x := y \rangle_i$, where y points to z . At the time of the assignment, z is reachable by definition,

⁷ The owner of an object may change over time.

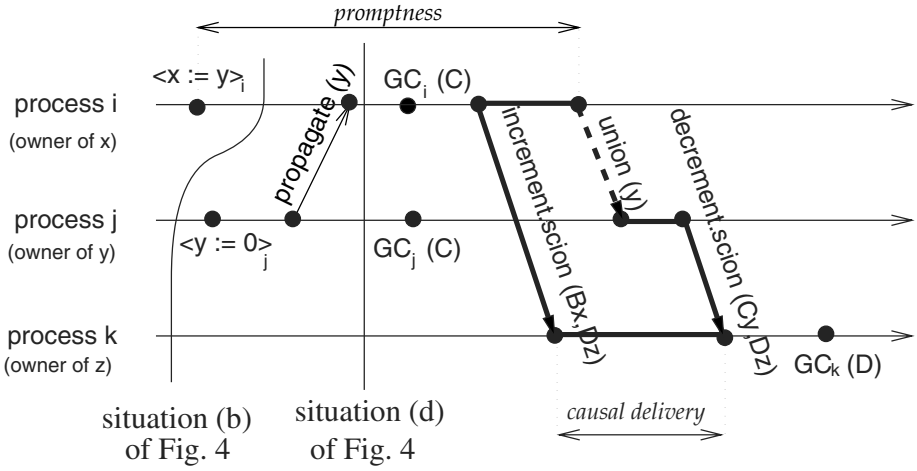


Fig. 5. Timeline showing the effect of some of the safety rules for the example of Figure 4. On site i, the sending of `increment.scion(Bx, Dz)` may be delayed at most until sending `union(y)`. Note the causal dependence (indicated by the thick lines) between the `increment.scion` and `decrement.scion` messages, carried by the `union` message.

and is guaranteed to remain reachable as long as y_i is not modified and remains reachable. It is not necessary for (a process managing) x to increment z 's reference count, as long as (some process managing) y does not decrement it.

Let us return to the example of Figure 4. At the time of $\langle x := y \rangle_i$, object z is reachable (from both replicas of y) and is protected by some scion, say `scion(Tt, -Zz)`; presumably, but not necessarily, $T = Y$ and $t = y$. As long as z 's scion has a non-zero count, it is safe to delay the increment of `scion(Bx, Dz)`. (Recall that it is the trace of bunch X which updates X 's set of stubs, which in turn causes the corresponding scion count to be adjusted.)

However, a problem remains with this approach. In the example, once situation (d) has been reached, it is possible that `decrement.scion(Cy, Dz)` reaches site k before `increment.scion(Bx, Dz)`; then z could be incorrectly reclaimed. To avoid this unsafe situation, it suffices to give precedence to `increment.scion` over `decrement.scion` and `union` messages. This is illustrated in Figure 5: the interval labeled *promptness*, shows how much the message `increment.scion(Bx, Dz)` can be delayed with respect to the moment when the corresponding assignment operation ($\langle x := y \rangle_i$) has been performed.

A set of further safety rules determine how late counting can be deferred, while still receiving messages in a safe order. We now state these rules, which will be justified in the following section.

Safety Condition II: Increment Before Decrement Rule. *Scanning an object causes the corresponding `increment.scion` messages to be sent immediately.*

Safety Condition III: Comprehensive Tracing Rule. *When process i sends a union or decrement.scion message, all replicas at i must have been scanned since most recently assigned.*

Safety Condition IV: Clean Propagation Rule. *When process i sends propagate(x), x_i has been scanned since it was most recently assigned.*

Safety Condition V: Causal Delivery Rule. *Garbage-collection messages (increment.scion, union and decrement.scion) are delivered in causal order.*

Rule II allows an object replica to be scanned at any time; scanning an object that contains a new pointer immediately sends an increment.scion message to the referent. It's important to mention that messages are asynchronous, so its actual transmission might take place later, provided messages are delivered in order, which is ensured by Rule V.⁸

Rule III ensures that union and decrement.scion messages are sent after increment.scion messages. In conjunction with Rule II, it ensures that, if an increment.scion message depends on a union or a decrement.scion, the latter will be sent before the former.

Rule IV ensures that when a process receives a new object via a propagate operation, any increment.scions corresponding to its new value have already been sent.

If delivery order is no better than FIFO, races can appear between increment.scion and decrement.scion messages. Rule V solves this problem. Note that coherence messages do not need causal delivery, thus limiting the cost.

Rules I through V are sufficient for the safe coexistence of replicated data and a hybrid garbage collector. They are independent of the coherence and tracing algorithms, and impose very few interactions between collection and coherence.

Justification of the Safety Rules. We now explain the preceding rules in more detail, by example.

Comprehensive Tracing Rule. This section justifies the Comprehensive Tracing Rule with an example of what happens if it is not enforced.

Consider Figures 4 and 5 in situation (d), i.e., after mutators have executed $\langle x := y \rangle_i$, $\langle y := 0 \rangle_j$, and y has propagated to site i . (Note that $\text{scion}(Bx, Dz)$ has not been created yet.) Suppose that $\text{trace}_i(C)$ runs and the Comprehensive Tracing Rule is not followed. The collector could send a union message to j (owner of y) indicating that $\text{stub}(Cy, Dz)$ has disappeared in process i , but not perform $\text{scan}_i(x)$ and $\text{increment.scion}(Bx, Dz)$ not sent. When j applies the Union Rule, it executes $\langle \text{send.decrement.scion}(Cy, Dz) \rangle_j$, causing $\text{scion}(Cy, Dz)$ to be deleted by k . Then, if $\text{trace}_k(D)$ runs, object z is reclaimed incorrectly.

⁸ In other words, sending a message only puts it on an ordered send queue.

The Comprehensive Tracing Rule prevents the above scenario because it forces x_i to be scanned before i sends the union message to j . Then, according to the Increment Before Decrement Rule, $\langle \text{send.increment.scion}(Bx, Dz) \rangle_i$ is performed before the union message is sent (and j applies the Union Rule and executes $\langle \text{send.decrement.scion}(Cy, Dz) \rangle_j$). Since we assumed causal delivery (Rule V) $\text{scion}(Bx, Dz)$ is created before $\text{scion}(Cy, Dz)$ is deleted. Consequently, z is not reclaimed by $\text{trace}_k(D)$.

Clean Propagation Rule. This section provides an example of what can happen when the Clean Propagation Rule is not enforced.

Consider Figures 4 and 5 after the mutator has executed $\langle x := y \rangle_i$ and before $\langle y := 0 \rangle_j$, i.e., in situation (b). Now, suppose that the coherence algorithm propagates the new value of x_i to some process w without first performing $\text{scan}_i(x)$. Then, the mutator executes $\langle x := 0 \rangle_i$. At this point, x no longer points to z , and the only scion that protects z is $\text{scion}(Cy, Dz)$. Suppose that both replicas of y are assigned in processes i and j to longer point to z either. By the collection algorithm, $\text{scion}(Cy, Dz)$ is deleted. Thus, z may be incorrectly reclaimed by $\text{trace}_k(D)$ (x_w still points to z).

The Clean Propagation Rule prevents the above scenario as it forces x_i to be scanned. Thus, by Rule II, $\langle \text{send.increment.scion}(Bx, Dz) \rangle_i$ is performed immediately, i.e., before x_i is propagated to site w .

Causal Delivery Rule. This section justifies the Causal Delivery Rule by yet another counter-example.

Consider Figure 4 in situation (d), i.e., after mutators have executed $\langle x := y \rangle_i$, $\langle y := 0 \rangle_j$, and y propagated to site i . (Note that $\text{scion}(Bx, Dz)$ has not been created yet.) Then, the collectors on sites i and j perform as follows: i executes $\langle \text{send.increment.scion}(Bx, Dz) \rangle_i$, whereas j executes $\langle \text{send.decrement.scion}(Cy, Dz) \rangle_j$. In an asynchronous system, the former could be delivered after the latter, causing z to be incorrectly reclaimed. In fact, there is a hidden causal relation through the shared variable y . In our algorithm, this causal relation is captured by the union message, as apparent in Figure 5. Thus, given the Causal Delivery Rule, there is at all times at least a scion protecting z from incorrect reclamation.

4.3 Discussion: A DGC Algorithm for Replicated Memory

What precedes focused on the interactions between garbage collection and replication (or caching), applied to a replicated (or cached) shared memory. We showed that both the tracing and the distributed counting garbage collector can execute independently of coherence. Garbage collection does not need coherent data, never causes coherence messages nor input/output, and it does not compete with applications' locks or working sets. However, coherence messages must at times be scanned before sending.

Our GC is a hybrid algorithm for a DSM. It combines tracing within a partition, with reference-counting across partition boundaries. Each process may

trace its own replicas, independently of other replicas. Counting (adjusting stubs and scions) at some process happens concurrently to other processes, and in the background with respect to the local mutator. In addition, counting is deferred and batched.

We presented five safety rules that guarantee the correctness of the distributed reference-counting algorithm. These safety rules are minimal and generally applicable:

- Union Rule: an object may be reclaimed only if it is unreachable from the union of all replicas (of the pointing objects);
- Increment before Decrement Rule: when an object is scanned, the corresponding `increment.scion` messages must be sent immediately;
- Comprehensive Tracing Rule: when a `union` or a `decrement.scion` message is sent, all replicas (on the sending site) must have been scanned since they were most recently assigned;
- Clean Propagation Rule: an object must be scanned before being propagated; and
- Causal Delivery Rule: GC messages must be delivered in causal order.

Measurements of our first (non-optimized) implementation [6] show that the cost of tracing is independent of the number of replicas, and that there is a clear performance benefit in delaying the counting.

Causal delivery, imposed by Rule V, is non-scalable in the general case; however, we do not consider this to be a serious problem in real implementations because causality can be ensured by taking advantage of the specific coherence protocols. For example, in our current implementation (supporting entry consistency) causal delivery is ensured by a mixture of piggy-backing and acknowledgments.

The Larchant algorithm is in use by the Esprit Project PerDiS [9], where it supports a large-scale cooperative engineering CAD application. This will enable us to measure and characterize the behaviour of real persistent applications, to fully study the performance of the distributed GC algorithm and to evaluate its completeness in a real-world environment. The PerDiS implementation is freely available at <http://www.perdis.esprit.ec.org/>.

5 Conclusion

We presented two recent advances in Distributed Garbage Collection in large-scale distributed computing systems. The first one is an algorithm for collecting cycles of garbage in a partitioned message-passing system. It is applicable in a clustered system (e.g., a network of LANs) communicating by RPC, as in Corba, Java RMI or DCOM. This algorithm has been implemented and proved correct. It is used in a prototype extension to ML supporting mobile objects, the JoCAML system [14].

Our second advance is an algorithm to collect garbage in a replicated or cached memory. It is applicable to distributed shared memories, object-oriented

databases, and persistent distributed stores. It has been proven safe and live, although by design for scalability, it is not complete. The algorithm provides the basis for the PerDiS platform developed in Esprit project 22.533 and used for cooperative engineering applications.

Acknowledgments

The research for both these algorithms was conducted within the Broadcast, Broadcast-WG, and PerDiS projects. We are grateful to Xavier Blondel for his improvements of Larchant and his research and implementation of DGC for the PerDiS platform. Many thanks to Steve Caughey and Ian Piumarta for their constructive comments.

References

- [1] Saleh E. Abdullahi and Graem A. Ringwood. Garbage collecting the internet: a survey of distributed garbage collection. *ACM Computing Surveys*, 30(3):330–373, September 1998. <http://www.acm.org/pubs/articles/journals/surveys/1998-30-3/p330-abdullahi/p330-abdullahi.pdf>.
- [2] Laurent Amsaleg, Olivier Gruber, and Michael Franklin. Efficient incremental garbage collection for workstation-server database systems. In *Proc. 21st Very Large Data Bases (VLDB) Int. Conf.*, Zürich (Switzerland), September 1995.
- [3] Henri E. Bal, Raoul Bhoedjang, Rutgwe Hofman, Criel Jacobs, Koen Langendoen, Tim Rühl, and M. Frans Kaashoek. Performance evaluation of the Orca shared-object system. *ACM Transactions on Computer Systems*, 16(1):1–40, feb 1998.
- [4] P.B. Bishop. Computer systems with a very large address space and garbage collection. Technical Report MIT/LCS/TR-178, Mass. Insitute of Technology, Cambridge MA (USA), 1977.
- [5] K. Mani Chandy and Leslie Lamport. Distributed snapshots: determining global states of distributed systems. *ACM Transactions on Computer Systems*, 3(1):63–75, February 1985.
- [6] Paulo Ferreira. *Larchant: ramasse-miettes dans une mémoire partagée répartie avec persistance par atteignabilité*. Thèse de doctorat, Université Paris 6, Pierre et Marie Curie, Paris (France), May 1996. http://www-sor.inria.fr/publi/ferreira_thesis96.html.
- [7] Paulo Ferreira and Marc Shapiro. Garbage collection and DSM consistency. In *Proc. of the First Symposium on Operating Systems Design and Implementation (OSDI)*, pages 229–241, Monterey CA (USA), November 1994. ACM. http://www-sor.inria.fr/publi/GC-DSM-CONSIS_OSDI94.html.
- [8] Paulo Ferreira and Marc Shapiro. Modelling a distributed cached store for garbage collection. In *12th Euro. Conf. on Object-Oriented Prog. (ECOOP)*, Brussels (Belgium), July 1998. http://www-sor.inria.fr/publi/MDCSGC_ecoop98.html.
- [9] Paulo Ferreira, Marc Shapiro, Xavier Blondel, Olivier Fambon, João Garcia, Sytse Kloosterman, Nicolas Richer, Marcus Roberts, Fadi Sandakly, George Coulouris, Jean Dollimore, Paulo Guedes, Daniel Hagimont, and Sacha Krakowiak. PerDiS: design, implementation, and use of a PERsistent DIstributed Store. Technical

- Report QMW TR 752, CSTB ILC/98-1392, INRIA RR 3525, INESC RT/5/98, QMW, CSTB, INRIA and INESC, October 1998. http://www-sor.inria.fr/publi/PDIUPDS_rr3525.html.
- [10] John Hughes. A distributed garbage collection algorithm. In Jean-Pierre Jouannaud, editor, *Functional Languages and Computer Architectures*, number 201 in Lecture Notes in Computer Science, pages 256–272, Nancy (France), September 1985. Springer-Verlag.
 - [11] Richard Jones and Rafael Lins. *Garbage Collection, Algorithms for Automatic Dynamic Memory Management*. Wiley, Chichester (GB), 1996. ISBN 0-471-94148-4.
 - [12] Leslie Lamport. Time, clocks, and the ordering of events in a distributed system. *Communications of the ACM*, 21(7):558–565, July 1978.
 - [13] Fabrice Le Fessant. The camlsspc system. <http://www-sor.inria.fr/projects/sspc/>, 1997.
 - [14] Fabrice Le Fessant. The jocaml system. Technical report, INRIA, 1998. <http://pauillac.inria.fr/join>.
 - [15] Fabrice Le Fessant, Ian Piumarta, and Marc Shapiro. An implementation for complete asynchronous distributed garbage collection. In *Proceedings of SIGPLAN'98 Conference on Programming Languages Design and Implementation*, ACM SIGPLAN Notices, Montreal, June 1998. ACM Press.
 - [16] T. Le Sergent and B. Berthomieu. Incremental multi-threaded garbage collection on virtually shared memory architectures. In *Proc. Int. Workshop on Memory Management*, number 637 in Lecture Notes in Computer Science, pages 179–199, Saint-Malo (France), September 1992. Springer-Verlag.
 - [17] Xavier Leroy. The objective-caml system software. Technical report, INRIA, 1996. <http://pauillac.inria.fr/ocaml>.
 - [18] Kai Li and Paul Hudak. Memory coherence in shared virtual memory systems. *ACM Transactions on Computer Systems*, 7(4):321–359, November 1989.
 - [19] Barbara Liskov, Mark Day, and Liuba Shrira. Distributed object management in Thor. In *Proc. Int. Workshop on Distributed Object Management*, pages 1–15, Edmonton (Canada), August 1992.
 - [20] U. Maheshwari and B. Liskov. Collecting distributed garbage cycles by back tracing. In *Principles of Distributed Computing*, Santa Barbara CA (USA), aug 1997. ACM.
 - [21] David Plainfossé and Marc Shapiro. A survey of distributed garbage collection techniques. In *Proc. Int. Workshop on Memory Management*, Kinross Scotland (UK), September 1995. http://www-sor.inria.fr/publi/SDGC_iwmm95.html.
 - [22] S. P. Rana. A distributed solution to the distributed termination problem. *Information Processing Letters*, 17:43–46, July 1983.
 - [23] Roger Sessions. *COM and DCOM: Microsoft's Vision for Distributed Objects*. Wiley, December 1998. ISBN 0-471-19381-X.
 - [24] Marc Shapiro, Peter Dickman, and David Plainfossé. SSP chains: Robust, distributed references supporting acyclic garbage collection. Rapport de Recherche 1799, Institut National de la Recherche en Informatique et Automatique, Rocquencourt (France), November 1992. http://www-sor.inria.fr/publi/SSPC_rr1799.html.
 - [25] Marc Shapiro and Paulo Ferreira. Larchant-RDOSS: a distributed shared persistent memory and its garbage collector. In J.-M. HéLary and M. Raynal, editors, *Workshop on Distributed Algorithms (WDAG)*, number 972 in Springer-Verlag LNCS, pages 198–214, Le Mont Saint-Michel (France), September 1995. http://www-sor.inria.fr/publi/LRDSPMGC_wdag95.html.

- [26] R. van Renesse, A. S. Tanenbaum, and A. Wilschut. The design of a high-performance file server. In *Proceedings of the 9th Int. Conf. on Distributed Computing Systems*, pages 22–27, Newport Beach CA (USA), June 1989. IEEE.
- [27] Paul R. Wilson. Uniprocessor garbage collection techniques. In *Proc. Int. Workshop on Memory Management*, number 637 in Lecture Notes in Computer Science, Saint-Malo (France), September 1992. Springer-Verlag. <ftp://ftp.cs.utexas.edu/pub/garbage/bigsurv.ps>.
- [28] Ann Wollrath, Roger Riggs, and Jim Waldo. A distributed object model for the java system. In *Conf. on Object-Oriented Technologies*, Toronto Ontario (Canada), 1996. Usenix.
- [29] V. Yong, J. Naughton, and J. Yu. Storage reclamation and reorganization in client-server persistent object stores. In *Proc. Data Engineering Int. Conf.*, pages 120–133, Houston TX (USA), February 1994.