

An On-the-Fly Reference-Counting Garbage Collector for Java

YOSSI LEVANONI, Microsoft Corporation

and

EREZ PETRANK, Technion - Israel Institute of Technology

Reference-counting is traditionally considered unsuitable for multi-processor systems. According to conventional wisdom, the update of reference slots and reference-counts requires atomic or synchronized operations. In this work we demonstrate this is not the case by presenting a novel reference-counting algorithm suitable for a multi-processor system that does not require any synchronized operation in its write barrier (not even a compare-and-swap type of synchronization). A second novelty of this algorithm is that it allows eliminating a large fraction of the reference-count updates, thus, drastically reducing the reference-counting traditional overhead. This paper includes a full proof of the algorithm showing that it is safe (does not reclaim live objects) and live (eventually reclaims all unreachable objects).

We have implemented our algorithm on Sun Microsystems' Java Virtual Machine 1.2.2 and ran it on a 4-way IBM Netfinity 8500R server with 550MHz Intel Pentium III Xeon and 2GB of physical memory. Our results show that the algorithm has an extremely low latency and throughput that is comparable to the stop-the-world mark and sweep algorithm used in the original JVM.

Categories and Subject Descriptors: D.1.5 [**Object-oriented Programming**]: Memory Management; D.3.3 [**Language Constructs and Features**]: Dynamic storage management; D.3.4 [**Processors**]: Memory management (garbage collection); D.4.2 [**Storage Management**]: Garbage Collection

General Terms: Languages, Algorithms.

Additional Key Words and Phrases: Programming languages, Memory management, Garbage collection, Reference-counting

1. INTRODUCTION

Automatic memory management is an important tool aiding in the fast development of large and reliable software systems. However, garbage collection does not come without a price—it has significant impact on the overall runtime performance. The amount of time it takes to handle allocation and reclamation of memory space may reach a noticeable percentage of the overall execution time for realistic programs. Thus, a clever design of efficient memory management and garbage collector is an important goal in today's technology.

Part of this work was presented at the *ACM Conference on Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA'01)*, October, 2001 [Levanoni and Petrank 2001]. Yossi Levanoni, Microsoft Corporation, One Microsoft Way, Redmond, WA 98052, USA. Email: ylevanon@microsoft.com.

Erez Petrank, Dept. of Computer Science, Technion - Israel Institute of Technology, Haifa 32000, Israel. Email: erez@cs.technion.ac.il.

This research was supported by the Coleman Cohen Academic Lecturship Fund and by the Technion V.P.R. Fund - Steiner Research Fund.

Most of this work was done while Yossi Levanoni was at the Dept. of Computer Science, Technion - Israel Institute of Technology, Haifa 32000, Israel.

1.1 Automatic memory management on a multiprocessor

In this work, we concentrate on garbage collection for multi-processor systems. Multi-processor platforms have become quite standard for server machines and are also beginning to gain popularity as high performance desktop machines. Many well studied garbage collection algorithms are not suitable for a multiprocessor. In particular, many collectors (among them the collector supplied with Javasoft's Java Virtual Machine) run on a single thread after all program threads (also known as *mutators*) have been stopped (the so-called *stop-the-world* concept). This technique is characterized by poor processor utilization, and hence hinders scalability.

In order to better utilize a multi-processor system, concurrent collectors have been presented and studied (see for example, [Baker 1978; Steele 1975; Dijkstra et al. 1978; Appel et al. 1988; DeTreville 1990; Boehm et al. 1991; Doligez and Leroy 1993; O'Toole and Nettles 1994; Doligez and Gonthier 1994; Printezis and Detlefs 2000]). A concurrent collector is a collector that does most of its collection work concurrently with the program without stopping the mutators. Most of the concurrent collectors need to stop all the mutators at some point during the collection, in order to initiate and/or finish the collection, but the time the mutators must be in a halt is short.

While representing a significant improvement over stop-the-world collectors, concurrent collectors also suffer from some shortcomings. Stopping all the mutators for the collection is an expensive operation by itself. Usually, the mutators cannot be stopped at any arbitrary point. Rather, they might be stopped only at *safe points* at which the collector can safely determine the reachability graph and properly reclaim unreachable objects. Thus, each mutator must wait until the last of all mutators cooperates and comes to a halt. On some systems, only mutator threads in an executing or runnable state must be stopped and therefore, the time to halt depends only on the number of processors, or on the number of non-blocked mutators, while on other systems all mutator threads must be handled until a full stop is obtained. This may hinder the scalability of the system, as the more mutators there are the more delay the system suffers. Furthermore, if the collector is not running in parallel (which is usually the case), then during the time the mutators are stopped, only one of the processors is utilized.

Therefore, it is advantageous to use *on-the-fly* collectors [Dijkstra et al. 1978; Doligez and Leroy 1993; Doligez and Gonthier 1994]. On-the-fly collectors never stop the mutators simultaneously. Instead, each mutator cooperates with the collector at its own pace through a mechanism called (soft) handshakes.

Another alternative for an adequate garbage collection on a multiprocessor is to perform the collection in parallel (see for example [Halstead 1985; Crammond 1988; Miller and Epstein 1990; Herlihy and Moss 1990; Endo et al. 1997; Flood et al. 2001; Kolodner and Petrank 2004]). We do not explore this avenue further in this work.

1.2 reference-counting on a multiprocessor

Reference-counting is a most intuitive method for automatic storage management. As such, systems using reference-counting were implemented more than forty years ago c.f. [Collins 1960]. The main idea is keeping for each object a count of the

number of references that reference the object. When this number becomes zero for an object o , it can be reclaimed. At that point, o may be added to the free list and the counters of all its predecessors (i.e., the objects that are referenced by the object o) are decremented, initiating perhaps more reclamations. It is possible to add the note to the free list and do the recursive deletion in a lazy fashion as illustrated in [Roth and Wise 1998]; this technique is likely to increase the locality of reference of the mutators.

Reference-counting seems very promising to future garbage collected systems. Especially with the spread of the 64-bit architectures and the increase in usage of very large heaps. Tracing collectors must traverse all live objects, and thus, the higher the utilization of the heap (i.e., the amount of live objects in the heap), the more work the collector must perform¹. Reference-counting is different. The amount of work is proportional to the amount of work done by the mutators between collections plus the amount of space that is actually reclaimed. But it does not depend on the space consumed by live objects in the heap.

The actual factors that determine which collector is better are complicated. They include heap occupancy, allocation rate, mutation rate, etc. However, some important program behaviors seem to match the advantages of reference-counting. Consider as an example a Web Server, servicing mostly static information out of local file systems and/or databases to external consumers. For a given request rate, it is beneficial for the server to have as much as possible of the requested data cached in memory. With tracing, as the size of this cache increases, the price paid for garbage collection increases, even though the load remains the same. With reference-counting, the price paid for garbage collection depends solely on the load and therefore allows utilizing more of the heap without incurring additional GC overhead.

The study and use of reference-counting on a multiprocessor has not been as extensive and thorough as the study of concurrent and parallel tracing collectors. The reason is that reference-counting has a seemingly inherent problem with respect to concurrency: the update of the reference-counts must be atomic since they are being updated by concurrent mutators. Furthermore, when updating a reference, a mutator must know the previous value of the reference slot being updated, i.e., the value it has overwritten. This value is hard to obtain in the presence of concurrent updates. If a mutator were to obtain an incorrect overwritten value, a mismatch is introduced into the bookkeeping of the reference-counts. Thus, the naive solution requires a lock on any update operation. More advanced solutions have recently reduced this overhead to a compare-and-swap operation for each reference update, which is still a time consuming overhead.

It is often claimed that even on a traditional uniprocessor client that runs single threaded programs reference-counting is more costly than tracing algorithms. One of the main reasons for that is the write-barrier overhead², which is high comparing

¹Other methods to ameliorate this problem, such as generational collectors were introduced. However, the full heap needs to be collected eventually. Subsequent work deals with generational collectors. See, for example, [Azatchi and Petrank 2003; Blackburn and McKinley 2003].

²The *write barrier* is the piece of code executed whenever a reference slot is modified. This code constitutes the overhead on each reference slot modification from a mutator.

to tracing algorithms that either require no write barrier or require a highly efficient (card-marking) write barrier when generations or concurrency are introduced.

1.3 This work

In this work, we present a new on-the-fly reference-counting garbage collector that ameliorates the two major problems of reference-counting mentioned above. First, the overhead of reference slots updates is reduced by introducing two paths in the write barrier: a fast path and a slow path. The slow path is as costly as the traditional overhead for reference slot update. However, the fast path is as fast as the write barrier used with concurrent tracing collectors. It turns out that for standard benchmarks the fast path is taken by the vast majority of reference slot updates, thus, reducing the cost of reference slot updates overhead to a level comparable with tracing collectors.

Second, the new algorithm has extremely fine synchronization. In particular, it avoids any synchronization in the write barrier. Thus, this algorithm is suitable for multithreaded programs on a multi-processor system. In addition, the new algorithm does not stop all mutators simultaneously. Instead, it suffices to stop one mutator at a time for cooperation with the collector.

Using the above two advantages, we derive a reference-counting collector with very short pauses, without paying a large cost in program throughput. We proceed with an overview on the novel ideas enabling these advantages.

The new algorithm, following Deutsch and Bobrow's *Deferred Reference-Counting* [Deutsch and Bobrow 1976], does not keep account of changes to local references (in stack and registers) since doing so is too expensive. Instead, it only tracks references in the heap (denoted *heap reference-count*). When garbage collection is required, the collector inspects all objects with heap reference-count zero. Those not referenced by the roots may be reclaimed. Our first observation is that in addition to local reference transactions, which are optimized away by the Deutsch-Bobrow technique, many more updates of the reference-counts are redundant and may be avoided. Consider a reference slot p that, between two garbage collections is assigned the values $o_0, o_1, o_2, \dots, o_n$ for objects o_0, \dots, o_n in the heap. All previous reference-counting collectors perform $2n$ updates of reference-counts accounting for these assignments: $RC(o_0)=RC(o_0)-1$, $RC(o_1)=RC(o_1)+1$, $RC(o_1)=RC(o_1)-1$, $RC(o_2)=RC(o_2)+1$, \dots , $RC(o_n)=RC(o_n)+1$. However, only two are required: $RC(o_0)=RC(o_0)-1$ and $RC(o_n)=RC(o_n)+1$. Building on this observation, it follows that in order to update all reference-counts of all objects before a garbage collection, it is enough to know which reference slots have been modified between the collections, and for each such slot, we must be able to determine its value in the previous garbage collection, and its current value.

In the new algorithm, we keep a record of all reference slots that have been modified. We also keep the “old” value that existed in the slot before it was first modified. It may seem problematic to obtain this value in a concurrent setting, and indeed, special care must be used to ensure that this value is properly registered. However, this is done without any synchronization operation. We denote the algorithm resulting from the discussion so far as *the snapshot algorithm*. This intermediate algorithm contains several of the novel ideas and is presented to distill these ideas clearly in this paper. The details of the snapshot algorithm are

presented in Section 4 below.

Next, we consider the collection itself. The naive implementation of the above approach is to stop all the mutators and read the values currently kept in all modified slots. This translates to taking a snapshot of the heap (or only a snapshot of the interesting slots in the heap). Such an approach does not allow full concurrency of the collector (it is not on-the-fly), although it is sound. In order to make the collector on-the-fly, we borrow ideas from the world of distributed computing. When taking a snapshot in a distributed environment, one does not stop all the computers over the distributed environment. Instead, one takes a snapshot of each computer at a time, but as the snapshots are being recorded, special care is taken to avoid confusion due to the non-instantaneous view. For example, all messages between computers are recorded as well. In our case, we will use an analogous solution. The analogy to a distributed system is that heap slots correspond to network nodes. The medium for communication between these nodes is the local state of mutators. Thus, as the algorithm obtains a non-instantaneous view of the heap slots (network nodes) it will also account for any information flowing into the slots from mutators' states (messages on-transit). We denote this non-instantaneous view of the heap *the sliding view*. The sliding view algorithm is described in Section 5.

1.4 Cycle collection

A disadvantage of reference-counting is that it does not collect all cycles. We have chosen to collect cycles with an on-the-fly mark-and-sweep collector. The Mark-and-sweep algorithm is seldom run to collect cycles and restore stuck reference-counts. (As in [Roth and Wise 1998; Wise 1993; Stoye et al. 1984; Chikayama and Kimura 1987], we use only two bits for the reference-count and thus, stuck counters are sometimes introduced. They are restored by the mark-and-sweep algorithm.)

We use a novel on-the-fly mark-and-sweep collector that we have designed especially for interoperation with the reference-counting algorithm. A description of this collector appears in [Azatchi et al. 2003]. It is quite natural to base a mark-and-sweep collector on a snapshot of the heap. The marking can be done on the snapshot view of the heap, and since unreachable objects remain unreachable, changes in the heap do not foil the collection of garbage. We adapt this basic idea to the sliding view notion, thus obtaining a tracing collector perfectly fitting our setting.

All measurements of throughput and latency in this paper are reported for the reference-count collector run most of the times and the on-the-fly mark-and-sweep run seldom. We describe our mark-and-sweep collector elsewhere [Azatchi et al. 2003].

1.5 Memory consistency

The algorithm presented in the paper requires a sequentially consistent memory. However, a few simple modifications can make the algorithm suitable for platforms that do not provide sequentially consistent memory. These modifications have negligible affect on the overall execution time. The Intel platform that we have used did not require these modifications. The modifications required and their cost are discussed in Section 6 below.

1.6 Implementation

We have implemented the algorithm on Sun's Reference Release 1.2.2 of the Java Virtual Machine and have run it on a 4-way IBM Netfinity 8500R server with a 550MHz Intel Pentium III Xeon processor and 2GB of physical memory. Two standard Java multithreaded benchmarks were used: SPECjbb2000 and the mtrt benchmark from SPECjvm98. These benchmarks are described in detail in SPEC's Web site [SPEC Benchmarks 2000]. As a sanity check, we have also run the collector on the client applications in SPECjvm98 on a uniprocessor.

1.7 Results

Section 9 reports performance measurements. It turns out that the algorithm has an extremely low latency, improving over the original JVM by two orders of magnitude. The measure reported is the one used by SPECjbb00: the maximum time it takes to complete a transaction. This is the same measure that was reported by Domani et al. [Domani et al. 2000] and it serves as an upper bound on the largest pause in execution time. The measured maximum transaction time was 16ms. The actual pause must be below that measure.

As for efficiency, on-the-fly collectors normally suffer a decrease in performance due to the added write barrier and synchronization [Printezis and Detlefs 2000; Ossia et al. 2002; Bacon et al. 2001]. Our collector shows excellent performance and does not fall much behind the compared stop-the-world tracing collector. In fact, for some benchmarks, the JVM with our reference-counting collector beats the original JVM by up to 10% improvement in the overall running times. This happens with the multithreaded mtrt benchmark. As for SPECjbb, if we allow a large maximum heap (which is the target of our collector), then our collector slightly improves over the running time of the original JVM. With smaller heaps the original JVM improves over ours by around 3% for SPECjbb.

1.8 Related work

The classic method of reference-counting, was first developed by Collins [Collins 1960]. It was later used in Small talk-80 [Goldberg and Robson 1983], the AWK [Aho et al. 1988] and Perl [Wall and Schwartz 1991] programs. Improvements on the original technique were suggested in several subsequent papers. Weizman [Weizenbaum 1963] ameliorated the delay introduced by recursive deletion. Deutsch and Bobrow [Deutsch and Bobrow 1976] eliminated the need for a write barrier on local references (in stack and registers). This method was later adapted for Modula-2+ [DeTreville 1990]. Further study on reducing work for local variables can be found in [Baker 1994] and [Park and Goldberg 1995]. An avenue for reducing redundant RC updates via static analysis is proposed in [Barth 1977]. Several works [Roth and Wise 1998; Wise 1993; Stoye et al. 1984; Chikayama and Kimura 1987] use a single bit for each reference-counter with a mechanism to handle overflows, under the assumption that most objects to be recovered are only singly referenced throughout their lives.

DeTreville [DeTreville 1990] describes a concurrent multiprocessor reference-counting collector for Modula-2+. This algorithm adapts Deutsch's and Bobrow's ideas of deferred reference-counting and transaction log for a multiprocessor system. However, the update operation is done inside a critical section that uses a single

central lock. So, no two updates can occur simultaneously in the system, placing a hard bound on its scalability.

Plakal and Fischer in [Plakal and Fischer 2000] propose a collection method based on reference counting for architectures that support explicit multithreading on the processor level. Their method requires co-routine type of cooperation between the mutator and a corresponding "shadow" collector thread and therefore is probably not suitable for stock SMPs, as SMP architectures do not support this kind of interaction in a natural and efficient manner.

Algorithms that perform garbage collection using a snapshot of the heap appear in [Furusou et al. 1991; Yuasa 1990]. In terms of synchronization requirements and characteristics our work is similar to that of Doligez, Leroy, and Gonthier [Doligez and Leroy 1993; Doligez and Gonthier 1994] in the sense that we use only fine synchronization, we never require a full halt of the system (the mutators are required to cooperate a few times per collection cycle). In our tracing algorithm we have used an object sweeping method similar to that presented in [Doligez and Leroy 1993; Doligez and Gonthier 1994].

Our garbage collection algorithm builds on a pragmatic concept of an atomic snapshot. The same concept is also a basic block in other garbage collectors, as well as in replication and checkpoint/restart algorithms. The pragmatism of this approach is dictated by the design of contemporary multiprocessor systems. There exists wide and profound theoretic research dealing with concurrent wait-free, shared memory algorithms that builds on the snapshot notion. The task in both research fields is common: to obtain a consistent picture of a system taken while it is conceptually frozen. Whereas much of the theoretical work assumes that threads are inherently faulty (hence these algorithms strive to be wait-free), however, the practical approach assumes that threads are not only inherently reliable, but also controllable by privileged roles (e.g., the operating system scheduler). Another difference is in the formulation of the problem—the theoretic community is more interested in the strain of the problem in which there are n threads, each with a single register to be shared among each other (a scatter-gather problem). The practical research, on the other hand, is concentrated on retrieving the values of all shared memory locations at a given instant, regardless of the number of running threads. For more reading on related theoretic research, the reader is referred to the work in [Riany et al. 1995] which provides an exciting attempt to bridge these two worlds as well as a comprehensive survey of previous atomic snapshot research.

1.8.1 *The work of Bacon et al.* Recent and novel work on reference-counting collectors, presented independently by Bacon et al [Bacon et al. 2001; Bacon and Rajan 2001], provides a solution to the entire garbage collection problem with a pure reference-counting approach, showing that this is feasible. Their contributions include a novel on-the-fly cycle collector, an improvement upon Deutsch and Bobrow's Deferred Reference-Counting algorithm that does not require use of Zero Count Tables, and an improvement over DeTreville's algorithm for concurrent collection by introducing epochs that eliminate the requirement for a single shared buffer for inc/dec operations.

Comparing the reported results of this algorithm to ours is not easy since the algorithms were implemented on different platforms. Three interesting points follow:

- Reducing synchronization.** A naive approach to multiprocessor reference-counting requires at least three compare-and-swaps in the write barrier. One for the update of the reference slot and two for the updates of the two reference-count. DeTreville [DeTreville 1990] has used a lock on each update to make sure that no two reference slot updates are executed concurrently. Bacon et al. [Bacon et al. 2001] made a significant step into exploiting multiprocessor concurrency by reducing the number of synchronizing operations to a single compare-and-swap. While significantly reducing the cost of synchronization, their write barrier still contains a compare-and-swap for each reference slot update. In this work, using the novel write barrier and sliding view ideas, we have managed to completely eliminate synchronization with reference slot updates. This major improvement is one of the more important contributions of this work.
- Throughput and latency:** It is not possible to compare the throughput of the two collectors since they have been run on different platforms and compared against different base JVM's. Both seem to be doing well with respect to gaining short pauses without sacrificing the throughput too much. With respect to pause times, the measured results provided by Bacon et al. are incomparable with ours. Bacon et al. used the Jikes JVM to measure the exact pause times. Unfortunately, we did not have the means to get such a measure, which is provided by the Jikes JVM. Instead, as in [Domani et al. 2000], we use the report output by the SPECjbb00 benchmark to indicate the maximum time it takes to complete a transaction. Both results show excellent latency with respect to previous reports of the same nature.
- Collecting cycles.** Finally, the two papers take different avenues for collecting cycles. Bacon and Rajan [Bacon and Rajan 2001] provide a novel approach for on-the-fly cycle detection. In contrast to their approach, we have chosen to develop an on-the-fly mark-and-sweep collector that exploits the sliding view mechanism and uses the same data structure as the reference-counting algorithm. This mark-and-sweep collector is run seldom in order to collect cycles and restore stuck reference-counts (see below).

1.9 Organization

In Section 2 a glossary of memory management terminology used in this paper is provided. In Section 3 definitions and terminology to be used in the rest of the paper are presented. In Section 4 the Snapshot Algorithm is presented. The sliding view algorithm is presented in Section 5. In Section 6 we discuss adaptation of the algorithm to platforms that do not provide sequentially consistent memory. In Section 7 implementation issues are discussed. In Section 8 we present the memory manager we used. In Section 9 we present performance results. In Sections 10 and 11 proofs are provided to the correctness of the presented collector. We conclude in Section 12.

2. A GLOSSARY

Here are some terms, common to the memory management literature and/or the Java Virtual Machine. This section is based on a similar section in [Hudson and Moss 2003] but is extended to match the terms appearing in this paper.

Class instance. See *Object*.

Collector, collector thread. A thread that performs garbage collection work (as opposed to executing application code).

Compare-and-swap (CAS). An atomic memory access operation, also known as compare exchange. $CAS(p, old, new)$ performs the following steps atomically. It compares the contents of location p with the value old . If the values are equal, it updates location p with the value new . The operation returns 1 if it installs the new value and 0 if it does not.

Concurrent garbage collection. Garbage collection methods in which the collector usually works concurrently with the mutators, but still requires hard handshakes (usually one or two per collection cycle—at the beginning and/or at the end of the cycle).

CPU; central processing unit. The hardware entity that executes instructions on behalf of threads. A computer system contains a fixed number of CPUs, all can read and write into a shared memory.

Cycles. Cycles of objects, such as circular lists. Reference-counting algorithms usually have difficulties reclaiming garbage objects which are linked in cycles.

Cycle collection. The ability (or lack thereof) to collect cycles of garbage.

Floating garbage. Garbage objects which are not being collected promptly after turning garbage. This subjective term usually denotes garbage that is generated during a collection cycle or garbage that a collection fails to reclaim even though eligible for collection.

Garbage; garbage objects; dead objects. Objects that are not reachable.

Global references. See *Roots*.

Handshake. A handshake is a synchronization mechanism between the collector and the mutators. In a handshake, the collector suspends the mutators and alters their local data structures while they're suspended. If all of the mutators are required to be suspended simultaneously, then this is considered a *Hard Handshake*. If only a single mutator is suspended at a time, the handshake is considered a *Soft Handshake*. Unless specified otherwise, in this paper, when the term handshake is used, it refer to a soft handshake. The collector may only suspend mutators when they reach safe points.

Heap reference-counting. A method of reference-counting in which only references from object slots to objects are counted. In this method, mutators' states are typically accounted for during collection cycles, and not during normal execution.

Latency. For a given benchmark, the average, median or maximal time it takes to complete a basic operation. In non-IO-bound systems, latency is highly correlated with pause times.

Limited field reference-counting. Reference-counting methods that allocate potentially overflowing reference-count fields to objects, in order to save space. Overflows are handled either by auxiliary lookup tables, which store the correct

count, or by fixing the count to a “stuck” value. Stuck counts may be restored by a tracing collector.

Local reference. A reference to an object that is present in a mutator’s local state.

Mark and sweep garbage collection. A garbage collection method that recursively traces and marks objects starting from the roots, then frees all non-marked objects.

Memory model. The set of guarantees that the hardware provides in terms of possible outcomes for programs involving concurrently executing threads. Sequential consistency is an example of a memory model. Traditional hardware supported sequential consistency while newer hardware designs reorder instructions (to increase CPU utilization) which results in a “weaker” memory model (that allows additional outcomes).

Memory synchronization operation. An operation that a CPU may issue in order to force perceived order of instructions throughout the system. Specifically, if a CPU issues instruction A, followed by perhaps additional instructions, followed by a memory synchronizing instructions, followed by perhaps some more instructions, followed by instruction B, then no CPU in the system will perceive instruction B as being executed before instruction A.

Mutators, mutator threads. Threads that perform application work. A mutator interacts with the collector when the mutator allocates objects and when it modifies reference slots. A mutator is also suspended sometimes by the collector during handshakes. As far as the algorithms in this paper are concerned, these are the relevant operations that a mutator may execute.

- Allocate a new object of a specified class. The result is a reference to the object placed in the mutator’s local state.
- Store a local reference into an object slot.
- Read an object slot into a local reference.

Mutator local state. A collection of object references and potentially other data structures that are immediately available to a mutator and may be used in the mutator operations outlined in the definition for *Mutator*, *mutator thread*. No other entity has access to a mutator’s local state, except for the collector: the collector may read and change the mutator’s state while the mutator is suspended.

Null reference. A distinguished reference value that refers to no object. In most systems, the null reference value is represented by a numerical 0.

Object; Class instance. An object is a collection of slots with a prescribed layout determined by the object’s class. Therefore it is also termed a class instance. In our model, objects exist only inside the heap. Objects contain two types of slots: scalars and reference. For the collector in this paper scalar slots are not relevant.

On-the-fly garbage collection. Garbage collection methods in which there is never more than a single mutator stopped. i.e., only soft-handshakes are used.

Additionally, it is usually required that this pause time be at most proportional to the size of the mutator's local state. (The sliding view algorithm, presented in this paper, is such a garbage collection algorithm.)

Parallel garbage collection. A method of garbage collection in which there is more than a single collector thread working on garbage collection. This is unrelated to whether the method is stop-the-world, concurrent or on-the-fly.

Pause time, maximal. The maximal duration a mutator has ever been stopped for by the collector. On-the-fly collectors display minimal pause times while stop-the-world collectors usually have longer pause times.

Protected code. Sections of code which are precluded for cooperation with the collector. i.e., sections of code that by definition do not contain safe points. The algorithms presented in this paper demand that the write barrier and object allocation procedures be protected.

Reachable. An object is reachable if a root refers to it, or another reachable object has a slot referencing it.

Read barrier. Operations performed when a mutator loads an object slot into its state, if such operations are required. The algorithms presented in this paper do not require a read barrier.

Reference-count. The number of references associated with a particular object by a particular reference-counting technique.

Reference-count field. A special system field associated with each object that contains the computed reference-count of the object. Depending on the reference-counting algorithm variant, the field may or may not contain the correct number of references to the object at all times. In addition, the physical placement of the field depends on both the reference-counting algorithm and the platform implementing the algorithm. For example, some algorithms or platforms will choose to place the field in-line with other instance fields of the objects while other algorithms and platforms will choose to place the field in a bitmap outside of the object.

Reference-counting. Garbage collection methods that determine reachability by counting the number of references to each object. The methods vary in whom is doing the counting (mutators vs. collector etc.) and what is being counted (references from heap slots solely or also references from mutators' states).

Reference slot. A slot within an object that refers to another object or contains the special value *null*.

Roots. A collection of variables that may contain references which are immediately accessible to at least one mutator. Thus, local references are roots. In most systems, there also exist static variables, intern tables, etc. which also provide immediate access to objects to mutators. To simplify the presentation of the algorithms in this paper, we assume that there exist a singleton object containing all such global variables. We then assume that each mutator has a local reference to this object. Thus, global references are modeled as object slots.³

³Our implementation, however, treats globals specially to reduce overhead.

Safe point. A point in the application code at which (a) accurate reference/non-reference type information is available for the mutator state (this is implementation dependent), and (b) no protected code is being executed. When a mutator is at a safe point, its GC-related information is consistent.

Sequential consistency. A memory model in which one can take the program traces of all threads and combine them into a sequential sequence in which (a) the order of instructions for each individual thread is preserved and (b) every load instruction returns the value stored into the location by the first store instruction into the location preceding the load in the sequential sequence. In other words, to all parties involved, the system appears to be working with a single CPU and without re-ordering of instructions.

Sliding view. A non-atomic picture of the shared memory of an SMP system. e.g., by copying each byte of the shared memory into disk over a period of time. A sliding view does not represent a coherent picture of the shared memory at any given moment. A sliding view provides a value for each location in the memory at some moment.

SMP; symmetrical multi-processor. A computer system possessing multiple CPUs and a shared memory. An SMP system is characterized by a memory model.

Snapshot. A copy of an SMP's system shared memory made while all CPUs are halted. i.e., a frozen and coherent picture of the shared memory.

Stop-the-world garbage collection. Garbage collection methods in which the collector and mutators seldom operate concurrently. These methods are usually characterized by long interruptions to user processing, yet are simple to implement and often exert lower overheads for the collection.

Throughput. For a given benchmark, the number of operations completed per time unit. One of two major measures of performance, together with latency.

Write barrier. Operations performed when a mutator stores an object reference from its local state into an object slot, if such operations are required. The algorithms presented in this paper require write barriers.

ZCT; zero count table. A set keeping track of objects whose heap reference-count is (or previously was) zero. These are candidates for collection, provided there is no local reference to them.

3. SYSTEM MODEL AND DEFINITIONS

For an introduction on garbage collection and memory management the reader is referred to [Jones and Lins 1996]. Slots in objects (in the heap) that hold references are denoted *heap reference slots* but most of the time they are just called slots. The reference-counts kept for each objects are *heap reference counts*. Namely, to count the number of references to an object, we sum over all slots that refer to this object in the heap. Mutators' local stacks and registers are not included in the count. As appropriate for Java, but also for other languages, we assume that all slots are initialized with a **null** reference when created. The reference-count associated with an object *o* are denoted *o.rc*.

Coordination between the collector and mutators. The algorithms presented in this paper assume that there is a single collector thread in the system and potentially multiple mutator threads.

It is assumed that the collector thread may *suspend* and subsequently *resume* mutators. When a mutator is suspended, the collector may inspect and change its local state with the effects taking place after the mutator is resumed. For example, a mutator may be stopped, and then buffers that it has written may be read, and new empty buffers may be placed instead of the filled buffers. The collector may then use the read buffers after reading them. When the mutator resumes, it may use the new empty buffers.

Garbage collectors in general, require that mutators are stopped for the collection at safe-points at which the collector can safely determine the reachability graph and properly reclaim unreachable objects. Classifying points in the JVM code into safe and unsafe depends on the specifics of the JVM. Since we run an on-the-fly collector, we do not need to stop all mutators together. However, when a mutator is suspended to read its local state (e.g., mark its roots), then we require that the mutator is at a safe point. In particular, we demand that it is not stopped during the execution of a reference slot update (including the write barrier code corresponding to the update) and it is not stopped during allocation of a new object via the `new` instruction.

4. THE SNAPSHOT ALGORITHM

For clarity of presentation, we start with an intermediate algorithm called *the snapshot algorithm*. Most of our novel ideas appear in it. In particular, a large fraction of the reference-counts updates (which are redundant) is saved and also a write barrier that does not require atomic operations or synchronization operations is introduced.

However, in this intermediate algorithm the mutators are stopped simultaneously for part of the collection. Therefore the pause lengths with this algorithm are not as good as with our final algorithm. The pauses imposed by this algorithm are not too long (the bottleneck is clearing a bitmap with dirty flags for all objects in the heap), but long enough to hinder scalability on a multiprocessor. In Section 5, this intermediate algorithm is extended to be made on-the-fly with very short pause.

The idea, as presented in Section 1.3, is based on computing differences between heap snapshots. The algorithm operates in cycles. A cycle begins with a collection and ends with another. A garbage-collection cycle is indexed by k and its actions are described next.

The first goal is to record all reference slots in the heap that have changed since the previous garbage collection (indexed $k - 1$), so that the corresponding reference-counts may be updated. The mutators do the recording with a write-barrier. In order to avoid recording slots again and again, a dirty flag is kept for each such slot. When a mutator updates a reference slot, it checks the dirty bit. If it is clear, the mutator records the slot's information in a local buffer, then sets the dirty bit. The recorded information is the address of the slot and its value before the current modification. Recording is done in a local buffer with no synchronization. Figure 1 provides the code for this operation.

```

Procedure Update(s: Slot, new: Object)
begin
1.   local old := read(s)
      // was s written to since the last cycle ?
2.   if  $\neg$ Dirty(s) then
      // ... no; keep a record of the old value.
3.   Bufferi[CurrPosi] :=  $\langle$ s,old $\rangle$ 
4.   CurrPosi := CurrPosi + 1
5.   Dirty(s) := true
6.   write(s, new)
end

```

Fig. 1. Mutator Code—Update Operation

```

Procedure New(size: Integer) : Object
begin
1.   Obtain an object o from the allocator,
      according to the specified size.
      // add o to the mutator local ZCT.
2.   Newi := Newi  $\cup$  {o}
3.   return o
end

```

Fig. 2. Mutator Code: for Allocation

When a collection begins, the collector starts by stopping all mutators and marking as local all objects referenced directly by the mutators' stack at the time of the pause. Next, it reads all the mutators' local buffers (in which modified slots are recorded), it clears all the dirty bits and it lets the mutators resume. After the mutators resume, the collector updates all the heap reference-counts to reflect their values at the time of the pause. Recall that the heap reference-count is the number of references to the object from other objects in the heap, but not from local variables on stack or registers. The algorithm for this update is presented and justified in the remainder of this section. However, if the heap reference-counts are properly updated, the collector may reclaim any object whose reference-count dropped to zero by this update as long as the object is not marked local. Such an object is not referenced from the roots, neither is referenced from other objects in the heap. As usual, the reference-counts of objects referenced by reclaimed objects are decremented and the reclamation proceeds recursively. A standard zero-count table (ZCT) [Deutsch and Bobrow 1976] keeps track of all objects whose reference-count drops to zero at any time. These objects are candidates for reclamation. Whenever an object is created (allocated) it has a zero heap reference count. Thus, all created objects are put in (a local) ZCT upon creation. The code for the create routine appears in Figure 2. New_i stands for the i^{th} mutator local ZCT—containing objects that have been allocated by T_i during the current cycle.

It remains to discuss updating the reference-counts according to all modified slots between collection $k - 1$ and k . As explained in Section 1.3, for each such slot s , we need to know the object O_1 that s referenced at the pause of collection

$k - 1$ and the object O_2 that s references at the pause of collection k . Once these values are known, the collector decrements the reference-count of O_1 and increments the reference-count of O_2 . When this operation is done for all modified slots, the reference-counts are updated and match the state of the heap at the k th collection pause.⁴

How to obtain the addresses of objects O_1 and O_2 is described next. We start with obtaining O_1 . If no race occurred when the slot s was first modified during this cycle, then the write barrier properly recorded the address of O_1 in the local buffer. It is the value that s held before that (first) modification and thus it is possible to obtain the address of O_1 and decrement its reference-count. But suppose a race did occur between two (or more) mutators trying to modify s concurrently. If one of the updating mutators sets the dirty flag of s before any other mutator reads the dirty flag, then only one mutator records this address and the recording will properly reflect the value of s at the $k - 1$ pause. Otherwise, more than one mutator finds the dirty bit clear. Looking at the write barrier code in Figure 1, each mutator starts by recording the old value of the slot, and only then it checks the dirty bit. On the other hand, the actual update of s occurs after the dirty bit is set. Under sequential consistency all threads observe the operations of other threads in program order. Thus, if a mutator detects a clear dirty bit it is guaranteed that the value it records is the value of s before any of the mutators has modified it. So while several mutators may record the slot s in their buffers, all of them must record the same (correct) information. To summarize, in case a race occurs, it is possible that several mutators record the slot s in their local buffers. However, all of them record the same correct value of s at the $k - 1$ st pause. When collecting the local buffers from all mutators multiple records of a slot are eliminated—they are all duplicates. We conclude that the address of object O_1 can be properly obtained by following the procedure just outlined.

It remains to obtain the address of O_2 , the object that s references at the pause of collection k . Note that at the time the collector tries to obtain this value the mutators are already running after the k th pause. The collector starts by reading the current value of s . It then reads s 's dirty flag. If the flag is clear then s has not been modified since the pause of collection k and we are done. If the dirty bit of s is set, then it has been modified. But if it has been modified, then the value of s at pause k is currently recorded in one of the mutators' local buffers. This value can be obtained by searching the local buffers of all mutators. It is not required to stop the mutators for peeking at their buffers. This slot must have a record somewhere and it will not be changed until the next $(k + 1)$ collection.

The collector operation is given in Figure 3. In Read-Current-State (Figure 4) the collector stops the s , takes their buffers, mark objects directly referenced from the roots as local, takes all local ZCT's (which include records of newly created objects), and clears all the dirty marks. The mutators are then resumed.

⁴It is normally a good practice to first do the increment and only then the decrement associated with a reference assignment. Since, in this work, reference counters are associated with snapshots, there is no danger in first decrementing due to the old value and only then incrementing due to the new value. The counts will remain non-negative and zero counts are added to the ZCT. Furthermore, decrementing first helps in avoiding stuck-counts when using small fields to represent

```

Procedure Collection-Cycle
begin
1.  Read-Current-State
2.  Update-Reference-Counters
3.  Read-Buffers
4.  Fix-Undetermined-Slots
5.  Reclaim-Garbage
end

```

Fig. 3. Collector Code

```

Procedure Read-Current-State
begin
1.  suspend all mutators
2.   $Hist_k := \emptyset$ 
3.   $Locals_k := \emptyset$ 
4.  for each mutator  $T_i$  do
    // copy buffer (without duplicates.)
5.   $Hist_k := Hist_k \cup Buffer_i[1 \dots CurrPos_i - 1]$ 
6.   $CurrPos_i := 1$ 
    // "mark" local references.
7.   $Locals_k := Locals_k \cup State_i$ 
    // copy and clear local ZCT.
8.   $ZCT_k := ZCT_k \cup New_i$ 
9.   $New_i := \emptyset$ 
10. Clear all dirty marks
11. resume mutators
end

```

Fig. 4. Collector Code—Procedure **Read-Current-State**

```

Procedure Update-Reference-Counters
begin
1.   $Undetermined_k := \emptyset$ 
2.  for each  $\langle s, v \rangle$  pair in  $Hist_k$  do
3.   $curr := \mathbf{read}(s)$ 
4.  if  $\neg Dirty(s)$  then
5.   $curr.rc := curr.rc + 1$ 
6.  else
7.   $Undetermined_k := Undetermined_k \cup \{s\}$ 
8.   $v.rc := v.rc - 1$ 
9.  if  $v.rc = 0 \wedge v \notin Locals_k$  then
10.  $ZCT_k := ZCT_k \cup \{v\}$ 

```

Fig. 5. Collector Code—Procedure **Update-Reference-Counters**


```

Procedure Read-Buffers
begin
1.    $Peek_k := \emptyset$ 
2.   for each mutator  $T_i$  do
3.     local  $ProbedPos := CurrPos_i$ 
       // copy buffer (without duplicates.)
4.      $Peek_k := Peek_k \cup$ 
        $Buffer_i[1 \dots ProbedPos - 1]$ 
end

```

Fig. 6. Collector Code—Procedure **Read-Buffers**

```

Procedure Fix-Undetermined-Slots
begin
1.   for each pair  $\langle s, v \rangle$  pair in  $Peek_k$ 
2.     if  $s \in Undetermined_k$  do
3.        $v.rc := v.rc + 1$ 
end

```

Fig. 7. Collector Code—Procedure **Fix-Undetermined-Slots**

The collector then proceeds to update the reference-counts (Figure 5). This happens after the s have been already resumed. As explained above, for each slot that is recorded in the mutators’ buffers (set $Hist_k$ in the algorithm) its contents are first read and then the dirty flag associated with it. If the dirty flag is off, then by the properties of the update barrier, and assuming sequential consistency, the read value for the slot is the prevailing one at the time the mutators were paused during the current cycle. If there is a “miss” then it is known that some mutator has modified s and has taken a record of its value in its current buffer. Such slots are put into the $Undetermined_k$ set.

The collector then executes procedure **Read-Buffers** (Figure 6) in which it obtains the set of slots that have already been modified by the mutators since they have been resumed, during the current cycle. Since the order of updating the local buffer and the dirty flag is (a) first write an entry into the local buffer (b) then increment the buffer pointer (denoted by $CurrPos_i$) (c) then set the dirty bit, it is guaranteed that for any slot in $Undetermined_k$, the collector is going to find a log entry in some mutator’s buffer. This is true because the collector has already noticed the raise of the dirty bit, hence it must also see the corresponding buffer entry and buffer pointer update.

In **Fix-Undetermined-Slots** (Figure 7) the collector resolves each of the slots in $Undetermined_k$ by looking them up against the $Peek_k$ set obtained in the previous step (procedure **Read-Buffers**). The corresponding reference-counters are incremented.

To summarize, the value of the “new” object referred to by s , denoted by O_2 is obtained in one of two ways. For most slots, the value of s is not modified by

reference-counts.

```

Procedure Reclaim-Garbage
begin
1.    $ZCT_{k+1} := \emptyset$ 
2.   for each object  $o \in ZCT_k$  do
3.     if  $o.rc > 0$  then
4.        $ZCT_k := ZCT_k - \{o\}$ 
5.     else if  $o.rc = 0 \wedge o \in Locals_k$  then
6.        $ZCT_k := ZCT_k - \{o\}$ 
7.        $ZCT_{k+1} := ZCT_{k+1} \cup \{o\}$ 
8.   for each object  $o \in ZCT_k$  do
9.     Collect( $o$ )
end

```

Fig. 8. Collector Code—Procedure **Reclaim-Garbage**

mutators during the short interval from the time the mutators resume and up until the collector tries to retrieve the value in procedure Update-Reference-Counters. All these slots may be read from the heap. Some slots however are fast-changing and will change during this time span. The collector will retrieve their values by peeking at the current mutators' buffers. At any rate, the update barrier provides a reliable mechanism for the collector to decide which route to take.

Finally, in procedure Reclaim-Garbage (Figure 8), the collector (recursively) reclaims all objects with zero reference-count that are not marked local. Every object with a zero reference-count but that is marked local is passed into the next cycle's ZCT. This is because it has a zero (heap) reference-count at the time that the mutators resume. In particular, it may cease to be locally referenced and need to be collected during the next cycle.

5. THE SLIDING VIEW ALGORITHM

The snapshot algorithm manages to execute a major part of the collection while the mutators run concurrently with the collector. A disadvantage of this algorithm is the halting of the mutators in the beginning of the collection. During this halt all mutators are stopped while the collector clears the dirty flags and receives the mutators' buffers and local ZCTs. This halt hinders both efficiency since only one processor executes the work and the rest are idle, and scalability since more mutators will cause more delays. While efficiency can be enhanced by parallelizing the flags' clearing phase, scalability calls for eliminating complete halts from the algorithm. This is indeed the case with our second algorithm, which avoids grinding halts completely.

A handshake [Doligez and Leroy 1993; Doligez and Gonthier 1994] is a synchronization mechanism in which each mutator stops at a time to perform some transaction with the collector. Our algorithm uses four handshakes. Thus, mutators are only suspended one at a time, and only for a short interval, its duration depends on the time to scan the roots.

In the snapshot algorithm there was a fixed point in time, namely when all mutators were stopped, for which the reference counts of all objects were computed. Thus, it was easy to claim that if an object has a zero heap reference-count at that

time, and it is not local at that time, then it can be reclaimed. By dispensing with the complete halting of mutators we no longer have this fixed point of time. Rather, there is a fuzziest picture of the system, formalized by the notion of a *sliding view* which is essentially a non-atomic picture of the heap. We show how sliding views can be used instead of atomic snapshots in order to devise a collection algorithm. This approach is similar to the way snapshots are taken in a distributed setting. Each mutator at a time will provide its view of the heap, and special care will be taken by the collector to make sure that while the information is gathered, modifications of the heap do not foil the collection.

5.1 Scans and sliding views

Pictorially, a scan σ and the corresponding sliding view V_σ can be thought of as the process of traversing the heap along with the advance of time. Each reference slot s in the heap is probed at time $\sigma(s)$; $V_\sigma(s)$ is set to the value of the probed reference slot. The *Asynchronous Reference-Count of o with respect to V_σ* is defined as follows.

DEFINITION 5.1. *For an object o and a sliding view V_σ we define the Asynchronous Reference-Count of o with respect to V_σ to be the number of slots in V_σ referring to o : $ARC(V_\sigma; o) \stackrel{\text{def}}{=} |\{s : V_\sigma(s) = o\}|$.*

Sliding views can be obtained incrementally with the benefit of not stopping all mutators simultaneously to compute the view. But in order to use this information safely we need to be careful. Trying to use the snapshot algorithm when we are only guaranteed that logging reflects some sliding view is bound to fail. For example, the only reference to object o may “move” from slot s_1 to slot s_2 , but a sliding view might miss the value of o in both s_1 (reading it after modification) and s_2 (reading it before modification).

These problems are avoided via a *snooping* mechanism. While the view is being read from the heap, the write-barrier marks any object that is assigned a new reference in the heap. These objects are marked as local, thus, preventing them from being collected in this collection cycle. (Recall that objects directly referenced by the roots are marked local to prevent collecting them even when they have a zero heap reference-count.) The snooped objects are left to be collected in the next cycle. The snooping mechanism, like Dijkstra’s write-barrier [Dijkstra et al. 1978], marks only objects that are reachable after the modification. Thus, it is less inclined to produce floating garbage than Yuasa’s write barrier [Yuasa 1990]. Assuming this *snooping* mechanism throughout the scan of the heap, we observe the following.

OBSERVATION 5.1. *If object o has $ARC(V_\sigma; o) = 0$, i.e., it is not referenced by any reference slot in the heap as reflected by the sliding view, and if object o is not referenced directly by the roots of the mutators after the scan was completed, and if object o has not been marked local by the snooping mechanism while the heap (and the roots) were being scanned, then at the time the heap scan is completed, object o is unreachable and may be reclaimed.*

PROOF. If the object is referenced by a heap slot in the end of the scan, then this slot has either been referencing this object when the scan of the heap read it,

or it has been written to that slot later. Both cases do not fall in the criteria of unreachable objects in the observation. Finally, if no reference is written into the heap while the roots are scanned, and there is no reference from the roots to this object, then it is unreachable. Here we rely on the fact that a mutator is stopped while reading its stack, so no reference may move while the mutator stack is being read; and furthermore, in Java, a reference cannot be moved from the stack of one mutator to another without being written to the heap. \square

Keeping this observation in mind, we are ready to present the sliding view algorithm. The description is broken into two parts. We first describe (in Section 5.2 below) how a sliding view of the heap may be used to reclaim unreachable objects. We call it a *generic* algorithm since it may use any mechanism for obtaining the sliding view. This is an extension of the ideas in the snapshot algorithm, still preserving the light write barrier. Then, in Section 5.3, a way to obtain a sliding view is presented.

5.2 Using sliding views to reclaim objects

Based on the above observation we present a generic reference counting garbage collection algorithm using sliding views.

Each mutator T_i has a flag, denoted $Snoop_i$ which signifies whether the collector is in the midst of constructing a sliding view. The mutators execute a write barrier when performing a heap slot update. After the store proper to the slot is performed, i.e., the reference to o is written into slot s , the mutator probes its $Snoop_i$ flag and, if the flag is set, the mutator marks o as *local*. This probing of the $Snoop_i$ flag and the subsequent marking is denoted *snooping*.

A collection cycle contains the following stages:

- (1) The collector raises the $Snoop_i$ flag of each mutator. This indicates to the mutators that they should start snooping.
- (2) The collector computes, using some mechanism, a scan σ and a corresponding sliding view, V_σ , concurrently with mutators' computations. (A possible implementation of this step appears in Section 5.3 below.)
- (3) Each mutator is then suspended (one at a time), its $Snoop_i$ flag is turned off and every object directly reachable from its roots is marked *local*. The mutator is then resumed.
- (4) Now, for each object o let $o.rc := ARC(V_\sigma; o)$.
- (5) Any object o that has $o.rc = 0$ and that was not marked *local* is reclaimed. For each reclaimed object, the reference-counts of all its descendants are decremented. If a descendant's reference-count is decremented to zero and the descendant is not marked local, then the descendant is reclaimed as well. This operation continues recursively until there are no objects that may be reclaimed.

A full proof of the sliding views algorithm is enclosed in Section 11 below. However, it is useful to state a few propositions which may help develop an intuition to why the above algorithm correctly reclaims unreachable objects.

PROPOSITION 5.2. *Let A be any allocated object and let T denote the time interval during which the sliding view is taken. If no reference to A is written to*

the heap or erased from the heap during T , then $ARC(V_\sigma; A)$ is exactly the heap reference-count of A at the end point of time interval T .

PROOF. If no reference to A changes during the interval T , during which the sliding view is taken, then reading the references to A at the end point of T yields the same information as reading these references any time during the time interval T . \square

PROPOSITION 5.3. *Let A be any allocated object and let T denote the time interval during which the sliding view is taken. If a reference to A is written to the heap during T , then A is not reclaimed during the current collection.*

PROOF. If a reference to A is written to the heap during T , then the snooping mechanism marks A local, and local objects are not collected. \square

PROPOSITION 5.4. *Let A be any allocated object and let T denote the time interval during which the sliding view is taken. If no reference to A is written to the heap during T , but some references to A are erased from the heap during T , then $ARC(V_\sigma; A)$ is greater or equal to the reference-count of A at the end point of the time interval T .*

PROOF. Let S be the set of reference slots pointing to A in the end of T . Since no reference to A is added, all slots in S refer to A at all points in the interval T , and so also when they are recorded by the sliding view. Thus, $ARC(V_\sigma; A) \geq |S|$, and $|S|$ (the size of the set S) is exactly the reference-count of A at the end point of time interval T . \square

PROPOSITION 5.5. *Let A be any allocated object and let T denote the time interval during which the sliding view is taken. Any object that has $ARC(V_\sigma; A) = 0$ and which is not snooped has heap reference-count 0 at the end of T .*

PROOF. If a reference to A has been written to the heap during T then it is snooped and we are done. If not, then either some heap references to A have been eliminated during T and then we are done by Proposition 5.4, or no heap references to A have been modified during T and then we are done by Proposition 5.2. \square

COROLLARY 5.6. *Let A be any allocated object and let T denote the time interval during which the sliding view is taken. If Object A is reclaimed during the current collection cycle, then A is unreachable when it is reclaimed.*

SKETCH OF PROOF. A is reclaimed if it has $ARC(V_\sigma; A) = 0$, it is not snooped, and it is not directly referenced by the roots when they are examined after T . Since $ARC(V_\sigma; A) = 0$ and A is not snooped, then by Proposition 5.5, A is not referenced by any heap reference slot at the end point of time interval T . Thus, to show that it is unreachable, it must be verified that A is not reachable directly from the roots. Suppose any mutator T holds a root reference to A at that time. If T writes this reference to the heap at any time until its local state is checked, then A will be snooped and will not be reclaimed, contradicting the assumption in the corollary. If T keeps the reference until its state is scanned, then A will be marked local and will not be reclaimed, again, contradicting the assumption in the claim. Finally, T may erase its local reference to A without writing it to the heap before its local

state is recorded. If that happens to all root slots that reference A then A is indeed unreachable when it is reclaimed.

It remains to check that the recursive deletion preserves the above properties. Indeed, all propositions above hold for the reduced heap obtained by removing any *unreachable* object and all its reference slots. They hold when the *ARC*'s are decremented for all descendants of the removed object. A full proof of this fact requires an induction on the order of removed objects. The induction hypothesis is that if a reference-count of the object A is decremented to zero during the reclamation process, then only references to it from unreachable objects existed during the end point of time interval T . \square

PROPOSITION 5.7. *Let A be any allocated object and let T denote the time interval during which the sliding view is taken. If A is not reachable before the snooping mechanism is started and if A is not part of a cyclic structure, then A is reclaimed during the current collection.*

PROOF. Since A is unreachable, then during the time interval T there are no modifications of reference slots referencing A or of reference slots referencing any of the objects from which A is reachable. By Proposition 5.2, the whole structure from which A is reachable will be assigned *ARC* that equals their heap reference-counts in the end of T . If there is no cycles in this structure, then the whole structure will be collected recursively by the reference-counting collector. \square

We have termed this algorithm “generic” since the mechanism for computing the sliding view in step (2) is unspecified. As in the snapshot algorithm, we do not intend to actually compute and record a full sliding view of the heap. Instead, we present an algorithm for updating the reference-counts for an implicitly defined sliding view of the heap. When the algorithm arrives at its collection phase, it holds for each object that $o.rc = ARC(V; o)$, where V is the sliding view that was constructed implicitly. Since we are not interested in the sliding view itself but rather on its manifestation through the *rc* fields, this implicit computation suffices for collection purposes.

5.3 Obtaining the sliding view

So far we have established the fact that garbage collection may take place even given only a sliding view of the heap and not a real snapshot. It is now shown how to obtain the sliding view. The algorithm of Section 5.2 above is used, but it is extended so that step (2), i.e., updating the reference-counts according to some sliding view, is done efficiently. The ideas originate from the snapshot algorithm described above. In particular, one does not need to record a full snapshot or sliding view of the heap. It is enough to record reference slot modifications and apply the corresponding reference-counts modifications. As in the snapshot algorithm, it is enough to know each slot's value in the previous sliding view and in the current one. The other reference slot modifications cancel out and are not recorded.

Four handshakes are used during the collection cycle to coordinate the collector with the mutators. In each handshake, each mutator is stopped, some coordination is run, and then the mutator resumes. The sliding view associated with a cycle spans from the beginning of the first handshake up to the end of the third handshake. The

“sampling” timing of each individual slot in the scan is determined by mutators’ logging regarding the slot. As dictated by the generic algorithm, the snooping flags are raised prior to the first handshake and are turned off at the fourth handshake, during which mutator local states are scanned.

Any reference slot which is modified between collection cycles is logged along with its value in the most recent sliding view so that the reference-counts may be updated to reflect the modifications of each modified slot. In contrast to the snapshot algorithm in which all dirty bits are cleared atomically, here we let the collector clear the dirty bits concurrently with the mutators. An important consequence of this relaxation is that slots may be inconsistently logged. For example, suppose the collector gets a mutator buffers in the first handshake and provides it with fresh new ones. At this time, the mutator modifies a slot s and logs its value in its new buffers. Next, the collector clears the dirty marks and just at that time a second mutator modifies s again and logs it in its local buffer with a different value. It turns out from our analysis that inconsistent logging of slots is only possible between responding to the first and responding to the third handshakes of a cycle. The important point here is that only one unique single value will be used for the sliding view. This value should be used to update the reference-counts at this collection (i.e., being the current value of s in the current sliding view) and also the same value should be used during the next collection when s ’s old value is sought. Just after the fourth handshake, the collector employs a consolidation mechanism to consolidate any inconsistently logged slot into one single value. No mutator would log a conflicting value after responding to the fourth handshake, hence no inconsistencies will be visible during the updates of the reference-counts or in the history recorded for the next cycle.

We proceed with the details of the sliding views reference-counting algorithm. In Section 5.4 the code executed by the mutator during reference slot modifications (the write barrier) and during object creation is presented. In Section 5.5 the collector algorithm is presented.

5.4 Mutator’s code

Mutators use the write barrier of the snapshot algorithm with the additional snooping and marking added after the store proper (see procedure **Update** in Figure 9). Object creation is the same as in the snapshot algorithm (see Figure 2 above).

5.5 Collector’s code

In this section, the main steps of the collection cycle are presented. The operation of each step is followed by the actual pseudo-code.

1. Signaling snooping. The collection starts with the collector raising the $Snoop_i$ flag of each mutator T_i , signaling to the mutators that it is about to start computing a sliding view.

2. Reading buffers (first handshake). Each mutator is stopped for a handshake. During the handshake mutators’ buffers are retrieved and then are cleared. (These are the same mutator buffers as in the snapshot algorithm.) The slots which are listed in the buffers are exactly those slots that have been changed since the last cycle. However, in the sliding view scenario this notion requires more care. The meaning of “changing” in this asynchronous setting is defined as follows. A slot

```

Procedure Update(s: Slot, new: Object)
begin
1.   Object old := read(s)
2.   if  $\neg$ Dirty(s) then
3.     Bufferi[CurrPosi] := ⟨s, old⟩
4.     CurrPosi := CurrPosi + 1
5.     Dirty(s) := true
6.   write( s, new)
7.   if Snoopi then
8.     Localsi := Localsi ∪ {new}
end

```

Fig. 9. Sliding View Algorithm: Update Operation

```

Procedure Initiate-Collection-Cycle
begin
1.   for each mutator Ti do
2.     Snoopi := true
3.   for each mutator Ti do
4.     suspend mutator Ti
       // copy (without duplicates).
5.     Histk := Histk ∪
       Bufferi[1 . . . CurrPosi - 1]
       // clear buffer.
6.     CurrPosi := 1
7.     resume Ti
end

```

Fig. 10. Sliding View Algorithm: Procedure **Initiate-Collection-Cycle**

is changed during cycle k if some mutator changed it after responding to the first handshake of cycle k and before responding to the first handshake of cycle $k + 1$. These are exactly the modifications that are obtained from the buffers at this stage. Steps (1) and (2) are carried out by procedure **Initiate-Collection-Cycle** (Figure 10).

3. Clearing. The dirty flags of the slots listed in the buffers are cleared. The slots are cleared *while the mutators are running*. This step is carried out by procedure **Clear-Dirty-Marks** (Figure 11). A critical difference between clearing the dirty bits while the mutators are all suspended and clearing them concurrently with the program, is that this step may clear dirty marks that have been concurrently set by the running mutators that have already responded to the first handshake. Since we want to keep these dirty bits set for the rest of this cycle, we will use the logging in the buffers (which currently contain all objects that have been marked dirty since the first handshake) to set these dirty bits on again in the next step.

4. Reinforcing dirty marks (second handshake). During the handshake the collector reads the contents of the mutators' buffers (which contain slots that were logged since the first handshake). The collector then *reinforces*, i.e., sets, the flags of the slots listed in the buffers.


```

Procedure Clear-Dirty-Marks
begin
1.   for each  $\langle s, o \rangle \in Hist_k$  do
2.      $Dirty(s) := \text{false}$ 
end

```

Fig. 11. Sliding View Algorithm: Procedure **Clear-Dirty-Marks**

```

Procedure Reinforce-Clearing-Conflict-Set
begin
1.    $ClearingConflictSet_k := \emptyset$ 
2.   for each mutator  $T_i$  do
3.     suspend mutator  $T_i$ 
4.      $ClearingConflictSet_k :=$ 
        $ClearingConflictSet_k \cup$ 
        $Buffer_i[1 \dots CurrPos_i - 1]$ 
5.     resume mutator  $T_i$ 
6.   for each  $s \in ClearingConflictSet_k$  do
7.      $Dirty(s) := \text{true}$ 
8.   for each mutator  $T_i$  do
9.     suspend mutator  $T_i$ 
10.  nop
11.  resume  $T_i$ 
end

```

Fig. 12. Sliding View Algorithm: Procedure **Reinforce-Clearing-Conflict-Set****5. Assuring reinforcement is visible to all mutators (third handshake).**

The third handshake is carried out. Each mutator is suspended and resumed with no further action. By the time all mutators resume, we know that they view correctly all dirty bits. Namely, a slot is dirty iff it was modified by a mutator that responded to the first handshake.

Steps (4) and (5) are executed by procedure **Reinforce-Clearing-Conflict-Set** (Figure 12).

6. Get-Local-States (fourth handshake). During the fourth handshake mutator local states are scanned and objects directly reachable from the roots are marked *local*. The snoop flags of the mutators are cleared. Mutators' buffers are retrieved once more and are *consolidated*. The last operation requires more explanation. As discussed in the beginning of Section 5.3 above, there may be some conflicting entries to a slot s in the buffers. The consolidation process takes care of leaving only one of these for future use of the collector.

Consolidating mutators' buffers amounts to the following. For any slot that appears in the mutators' buffers accumulated between the first and fourth handshakes, pick *any* occurrence of the slot and copy it to a digested consistent history. All other occurrences of the slot are discarded.

Intuitively, the sliding view value of any slot s that goes through the consolidation process is exactly the consolidated value. The reference count of the object referenced by this value will be incremented in this cycle. Thus, it is important

```

Procedure Get-Local-States
begin
1.   local  $Temp := \emptyset$ 
2.    $ColLocals_k := \emptyset$ 
3.   for each mutator  $T_i$  do
4.     suspend mutator  $T_i$ 
5.      $Snoop_i := \mathbf{false}$ 
      // copy and clear snooped objects set
6.      $ColLocals_k := ColLocals_k \cup Locals_i$ 
7.      $Locals_i := \emptyset$ 
      // copy mutator local state and ZCT.
8.      $ColLocals_k := ColLocals_k \cup State_i$ 
9.      $ZCT_k := ZCT_k \cup New_i$ 
10.     $New_i := \emptyset$ 
      // copy local buffer for consolidation.
11.     $Temp := Temp \cup$ 
       $Buffer_i[1 \dots CurrPos_i - 1]$ 
      // clear local buffer.
12.     $CurrPos_i := 1$ 
13.    resume mutator  $T_i$ 
      // consolidate  $Temp$  into  $Hist_{k+1}$ .
14.     $Hist_{k+1} := \emptyset$ 
15.    local  $Handled := \emptyset$ 
16.    for each  $\langle s, v \rangle \in Temp$ 
17.      if  $s \notin Handled$  then
18.         $Handled := Handled \cup \{s\}$ 
19.         $Hist_{k+1} := Hist_{k+1} \cup \{\langle s, v \rangle\}$ 
end

```

Fig. 13. Sliding View Algorithm: Procedure **Get-Local-States**

to make sure that its reference-count will be decremented in the next cycle, if s is modified again. Therefore, the digested history output by the consolidation mechanism replaces the accumulated mutators' buffers. Namely, the history for the next cycle is comprised of the digested history of mutators' logging between the first and fourth handshakes of the current cycle, unified with mutators' buffers representing updates that will occur after the fourth handshake of the current cycle but before the first handshake of the next cycle. The above operations are carried out by procedure **Get-Local-States** of Figure 13.

7. Updating. The collector proceeds to adjust rc fields due to differences between the sliding views of the previous and current cycle. This is done exactly as in the snapshot algorithm (see Figure 5). The collector will fail to determine the “current” value of all slots that were modified (i.e., are dirty). These slots will be treated later and are now marked as *undetermined*.

8. Gathering information on undetermined slots. The collector asynchronously reads mutators' buffers (using the snapshot algorithm procedure **Read-Buffers** of Figure 6). Then, in procedure **Merge-Fix-Sets** (Figure 14) it unifies the set of read pairs with the digested history computed in the consolidation step. The set of undetermined slots is a subset of the slots appearing in the unified set so

```

Procedure Merge-Fix-Sets
begin
1.    $Peek_k := Peek_k \cup Hist_{k+1}$ 
end

```

Fig. 14. Sliding View Algorithm: Procedure **Merge-Fix-Sets**

```

Procedure Collect(o: Object)
begin
1.   local DeferCollection := false
2.   foreach slot s in o do
3.     if Dirty(s) then
4.       DeferCollection := true
5.     else
6.       val := read(s)
7.       val.rc := val.rc - 1
8.       write(s, null)
9.       if val.rc = 0 then
10.        if val  $\notin$  ColLocalsk then
11.          Collect(val)
12.        else
13.           $ZCT_{k+1} := ZCT_{k+1} \cup \{val\}$ 
14.        if  $\neg DeferCollection$  then
15.          return o to the general purpose allocator.
16.        else
17.           $ZCT_{k+1} := ZCT_{k+1} \cup \{o\}$ 
end

```

Fig. 15. Sliding View Algorithm: Procedure **Collect**

the collector may now proceed to look up the values of these undetermined slots.

9. Incrementing *rc* fields of objects referenced by undetermined slots. In procedure **Fix-Undetermined-Slots** (same one from the snapshot algorithm—Figure 7) any undetermined slot is looked up in the unified set and the *rc* field of the associated object is incremented.

10. Reclamation. Reclamation generally proceeds as in the previous algorithm, i.e., recursively freeing any object with zero *rc* field which is not marked *local*. However, we do not reclaim objects whose slots appear in the digested history. These are objects which were modified since the cycle commenced but became garbage before it ended. Since the buffers will be used in the next collection, the reclamation of such objects is deferred to the next cycle. Reclamation is carried out using the procedures **Reclaim-Garbage** (same one as in the snapshot algorithm—Figure 8 and **Collect** (Figure 15).

5.6 Some words on correctness

A full proof of the collector appears in Sections 10-11 below. However, some preliminary discussion may help understand the algorithm and its properties. This algorithm is an extension of the generic algorithm presented in Section 5.2. It is

shown there that it is possible to collect garbage safely using a sliding view. It remains to show that the full algorithm indeed collects according to some sliding view of the heap. As in the snapshot algorithm, the sliding view is not explicitly written or recorded. Instead, only modifications in the view are recorded and they are enough to compute the updated asynchronous reference-counts of all objects. The first step is to define the sliding view that is implicitly used by the algorithm. Then, it is claimed that the algorithm indeed computes the asynchronous counts appropriately according to this sliding view.

5.6.1 *The sliding view associated with a cycle.* To define a sliding view, one must define a time in which each slot is probed. Recall (as in Section 5.1 above) that $\sigma(s)$ is the time at which the slot s is scanned, and $V_\sigma(s)$ is the value probed during the scan. Consider any memory word s . We define the scan time of s as follows.

- **Case 1:** If s does not appear in the buffers that are read during the first handshake of collection k , then s does not cause any reference-count updates during the current collection. We choose to fix its scan time to the beginning of the first handshake.
- **Case 2:** If s does appear in $Hist_k$, we split the discussion into two cases.
 - **Case 2a:** The slot s is logged by some mutator T_i after T_i responds to the first handshake and before it responds to the third handshake. In this case, the value of s for $Hist_{k+1}$ is later set by the consolidation mechanism in procedure **Get-Local-States**. In this case, the scan time is defined according to the consolidated value for s . Specifically, there is a value that survives the consolidation process for this slot, and there is a mutator that reads this value and stores it into its buffer. The scan time is fixed to be the time this mutator read the said value.
 - **Case 2b:** Otherwise, s is not logged by any mutator between the first and third handshakes and then we choose the scan time to be the end of the second handshake.

Intuitively, Case (1) is simpler since the collector does not clear the dirty bits of slots that are not in $Hist_k$. In this case, the write barrier works similarly to the snapshot algorithm. More care will be needed when $s \in Hist_k$ and the dirty bit is cleared during the run. In Case (2) we distinguish between slots that may have been logged with two different values because of the dirty bit clearance and slots that must have been logged with a unique value. The following simple assertions about the sliding view are presented without proof.

PROPOSITION 5.8. *The scan of the sliding view of cycle k start at or after the beginning of the first handshake and it ends at or before the end of the third handshake.*

COROLLARY 5.9. *The snooping mechanism is active throughout the scan of the sliding view.*

5.6.2 *The collector uses the sliding view.* The per-cycle sliding view is computed implicitly by the collector and mutators (bearing similarity to the conceptual snapshot taken at collection k by the snapshot algorithm which is never explicitly

computed.) We say that a slot is modified if it is modified by some mutator after responding to the first handshake of the previous cycle (cycle $k - 1$) and before responding to the first handshake of the current cycle (cycle k). If a slot is not modified, then it does not cause any reference-count modification to its referent. We first claim that this is correct. In other words, the value of such a slot at the current sliding view is equal to its value in the previous sliding view.

PROPOSITION 5.10. *If $s \notin Hist_k$ then s is not modified between the time s is scanned for the sliding view of collection $k - 1$ and the time s is scanned for the sliding view of collection k .*

PROOF. We know that slot s is not modified by any mutator after responding to the first handshake of the previous cycle and before responding to the first handshake of the current cycle. We show that the value of s at the beginning of the first handshake of cycle k is the value of s in both sliding views. Starting with the cycle k , the assertion holds by definition since the scan of s is defined to be in the beginning of the first handshake. To see that the same value is also read by the scan of cycle $k - 1$, we divide the analysis into two cases. If s has been modified in the cycle before, i.e., $s \in Hist_{k-1}$, then the scan of s in cycle k_1 happens after the first handshake of that cycle ends. This is true both in case (2a) and case (2b). Now, using the fact that s is known not to change after the first handshake of cycle k_1 , we are done. The remaining case is when s is also not in $Hist_{k-1}$. But then it has not been modified for two consecutive cycles and scanning it in the beginning of the first handshake of cycle $k - 1$ (as dictated by Case 1) yields the same value as in any other point in time while it remains unchanged. \square

We conclude that if s is not in $Hist_k$ then no updates are necessary according to s on the asynchronous reference-counts. This is very much like the computation of the reference-counting of the snapshot algorithm. We now turn to slots that are modified.

PROPOSITION 5.11. *If $s \in Hist_k$, then the collector properly updates the reference-counting changes caused by modifying s between the sliding view of the previous cycle to the sliding view of the current cycle.*

SKETCH OF PROOF. The collector reads the value of s in the previous cycle from the buffers, i.e., through $Hist_k$. We need to show that this recorded value equals the value of s in the sliding view of cycle $k - 1$ as defined above. Then, the collector obtains the current value from the heap or from the new buffers. We need to show that the value obtained matches the value of the sliding view of cycle k . We start with the latter.

If the collector reads s in Procedure **Update-Reference-Counts** and finds it not dirty, then s has not been modified by any mutator after responding to the first handshake. This is true because all dirty bits are reinforced before the third handshake and the collector reads these values after the fourth handshake. Thus, in this case, the same value is fixed throughout the scan time and in particular its value matches the one defined according to Case (2b).

If the collector finds s dirty, then one possibility is that it falls into Case (2a), i.e., it has been modified by a mutator between responding to the first and the third handshake. In this case, the consolidation mechanism fixes the value of this slot.

The consolidated value is the one read by the collector, and it is by definition the value of the sliding view. So this case is fine. Otherwise, the value of s has been modified by a mutator after responding to the third handshake. In this case, the value of s is not changed between the end of the first handshake and the beginning of the third handshake. In addition, the write barrier is guaranteed to record a unique value to the buffers, which, similarly to the proof of the snapshot algorithm, is the value existing in this slot before the third handshake begun. Indeed, the definition of the scan for s in this case, determines the end of the second handshake to be the scan time, thus fixing the value of slot s in the current sliding view to exactly the value obtained from the buffers in Procedure **Read-Buffers**.

It remains to check that the value of s found in $Hist_k$ matches the value of s in the scan of the cycle $k - 1$. Here, again, we need to partition the analysis into two cases according to whether $s \in Hist_{k-1}$, i.e., whether it was modified between cycles $k - 2$ and $k - 1$.

Suppose $s \in Hist_{k-1}$. If it was modified between the first and third handshakes of the cycle $k - 1$, then its scan value is determined according to the consolidated value of cycle $k - 1$. The same value exactly is inserted into $Hist_k$ by the consolidation mechanism and we are done. If s was not modified during the first and third handshake of cycle $k - 1$ (yet, s has been logged into $Hist_k$ later), then the write barrier properties allow only one value to appear in the buffers. This value is the value that s held between the end of the first handshake in cycle $k - 1$ and the beginning of the third handshake in that cycle. By case (2b) this value is scanned for the sliding view of cycle $k - 1$ during this interval (by the end of the second handshake). Thus, the sliding view value is determined to be exactly the value that appears with s in $Hist_k$ and we are done.

If $s \notin Hist_{k-1}$ then the dirty bit is not cleared during collection $k - 1$. In this case, the write barrier behaves exactly as in the snapshot algorithm and we are guaranteed to see in $Hist_k$ the value of the slot s as it was between cycles $k - 2$ and $k - 1$, thus, also in the beginning of the first handshake in cycle $k - 1$ when it is scanned for the sliding view of cycle $k - 1$ and we are done. \square

The propositions above show that the collector is properly updating the asynchronous reference-counts according to the sliding views defined. By the properties of the generic algorithm described in Section 5.2, we get that the collector safely collects garbage.

6. MEMORY COHERENCE

As mentioned in the introduction, a few modifications can make the algorithm suitable for platforms that do not guarantee sequential memory consistency. We list these modifications here and discuss their cost. Note that the collector is single threaded. Therefore, the concern is about undesirable interaction between two concurrent mutator threads and between the collector thread and a mutator thread.

The algorithm relies on the following order of instructions in the write barrier, in the (more interesting) case in which the dirty flag is found off and the snoop flag is found on:

- (1) Read the value v of slot s
- (2) Read the value of the dirty flag (returns: clean)

- (3) Write a record $\langle s, v \rangle$ to the local buffer (may require multiple memory accesses)
- (4) Update the buffer pointer
- (5) Dirty the flag
- (6) Write $newv$ into s
- (7) Read snoop flag (returns: on)
- (8) Snoop $newv$

This order of instructions has the following significance in the algorithm:

Dependency 1: Had any of the instructions (1), (2), (5) or (6) been reordered, the value associated with s in the buffer might have been incorrect. In addition, this might have foiled the collector’s effort to determine the slot’s value. The interdependent pieces of code for this instructions order are (a) this code itself (i.e., mutators racing to update a slot) (b) in the collector, trying to determine a slot. The collector actually executes a prefix of the above code (instructions (1), (2) and (3)) in procedure Update-Reference-Counters (see Figure 5) and has the same ordering requirements.

Dependency 2: Writing into the log (in instructions (3) and (4)) should occur strictly before dirtying the flag in (5) or else the collector may not be able to retrieve the values of undetermined slots in procedure **Read-Buffers** (Figure 6). Instructions (3) and (4) themselves may not be reordered because otherwise the collector may read half-updated log entries.

Dependency 3: Reading of the snoop flag in (7) must occur after the store proper in (6) or else some slots may be updated after the collector has started a collection cycle, without being snooped. This breaks the safety of the algorithm.

Solving dependencies (2) and (3) can be done using two lightweight handshakes. One letting the mutators acknowledge they have sensed the raising of the snoop flag before the collector starts the first handshake (for dependency 3), and the other occurring during (or prior to) the procedure Read-Buffers to make sure that the mutators have checkpointed the value of the buffer pointer (for dependency 2). In this lightweight handshake the collector updates a “stage” flag and the mutators follow suite on their earliest convenience, after having issued a memory synchronizing instruction. From the collector side, after a mutator signals it has cooperated, it is guaranteed that it has picked up the value of the snoop flag or that it has made available its most recent log buffer updates to the collector. Adding two lightweight handshakes to a collection cycle causes a negligible cost to the throughput.

To solve dependency (1), we use a mixture of techniques. First we handle the mutator side of this dependency. In our implementation, the write barrier begins with an optimistic test of the dirty bit. Only if the dirty bit is off, does the mutator proceed to execute the full write barrier as described in the sliding view algorithm. If the dirty bit is set, then the mutator may proceed immediately to the actual reference slot update and to the snooping check. It turns out from our measurements that for typical benchmarks most reference slot updates are executed on newly created objects. In our implementation, objects are marked dirty, and are logged, as soon as they are created, with a special compressed record denoting that they point to no object. Thus, the “heavy” write barrier path is not executed when an object is initialized (soon after it is created.) Since most objects are short

```

Procedure Optimistic-Update(s: Slot, new: Object)
begin
1.   if  $\neg$ Dirty(s) then
2.     Object old := read(s)
3.     Synchronization barrier
4.     if  $\neg$ Dirty(s) then
5.       Bufferi[CurrPosi] :=  $\langle$ s, old $\rangle$ 
6.       CurrPosi := CurrPosi + 1
7.       Dirty(s) := true
8.       Synchronization barrier
9.     write( s, new)
10.  if Snoopi then
11.    Localsi := Localsi  $\cup$  {new}
12.  end

```

Fig. 16. Sliding View Algorithm: Optimistic Update Operation

lived, this eliminates the need to call the heavy path of the write barrier on the striking majority of cases. Our measurements show that for the javac benchmark this happens less than once in a hundred, and for the SPECjbb benchmark and all the other SPECjvm98 benchmarks this happens less than once in a thousand. We will incorporate synchronization only in the heavy path of the write barrier, i.e., when the object is found not dirty. thus, the vast majority of the reference slot updates require no cost for handling memory coherence. In order to enforces the required instructions order, memory synchronizing instructions are placed between instructions (2) and (4) and between instructions (7) and (9) of the update barrier. See Figure 16 for the revised code. This indeed renders it heavyweight but as said, this code is seldom executed.

On the collector side, when trying to determine a slot, the same technique cannot be used efficiently. This is because mutators usually update new objects (for initialization), which are initially dirty, whereas the collector tries to determine the values of slots of old objects, that have survived at least one collection cycle and whose dirty bit has therefore been reset at the initial stages of the current cycle. Hence most of the dirty bits the collector probes are set to false.

Rather, a batching technique should be used: instead of trying to determined the value of a single slot at a time, the collector tries to determine the values of a batch of N slots at a time. These are the steps the collector follows:

- (1) Read N slots from the history buffer for the cycle.
- (2) Read the values of the slots, storing them in a local array.
- (3) Place a memory synchronizing instruction.
- (4) For each slot, read the value of the corresponding dirty flag and act appropriately, using the recorded value for the slot in the local array.

By choosing N to be large enough, the price of memory synchronization can be reduced to a very low level. Batching may or may not reduce data locality. In fact, it may even help increase locality as similar types of access are grouped together (slot retrieval vs. reference-count updates).

We have not implemented these modifications (except for measuring the frequency of entering the “heavyweight” section of the update barrier), since they were not required by the Intel platform we had used for testing.

7. IMPLEMENTATION ISSUES

In this section we shift from the abstract treatment of the dirty flags and the log buffers and suggest concrete implementations for these data structures.

7.1 Global roots

For convenience, we assume in the exposition of the algorithms and their proofs that there are no global roots. Instead, we model global roots as members of a distinguished heap object which is reachable from the local state of every mutator. This is enough to make the algorithm correct.

In practice, we give a special treatment to global roots by using a designated write barrier with them. Global roots are not assigned a dirty bit because they are scanned on each cycle regardless of whether they have been modified or not. Therefore the designated write barrier does not mark slots dirty when they change. However, it does invoke the snooping mechanism. Thus, each new reference written to a global root while the snoop flag is set, makes the referent snooped. Such an object cannot be collected during the current cycle. Finally, when local states are checked during the fourth handshake, objects reachable from global roots are also marked local.

7.2 Dirty Flags

This section explains how one dirty bit can protect several heap slots, and why that bit became an instance variable for each object.

We start with explaining how a flag can serve an indicator to a change in *any* of the slots within a fixed chunk of memory. The ideas are similar to those that arise in the context of tracking inter-generational references in a generational collector that uses *card marking* [Sobalvarro 1988].

In order to let a single flag signify a change in a chunk of memory we let the write barrier take the following form. When the value v is stored into the slot s , the flag for the chunk of memory containing s is probed. If the flag is set then the mutator may proceed directly to the store operation. Otherwise, the following cautious action is taken.

- A replica of the slots that reside inside the chunk is created and stored locally.
- the flag for the chunk is then probed again. If it is now turned on, we proceed to the store operation, without using the local replica.
- otherwise, i.e., the dirty flag is still clear, we copy the replica just prepared to the log buffer, raise the flag and only then execute the store.

The collector code for determining a slot is changed accordingly. The collector tries to determine the value of an entire chunk instead of a single slot. Using one dirty bit for a chunk of memory is characterized by a decreased memory consumption (less dirty bits) yet by spurious work imposed on the mutator and collector that have to process slots which haven’t really changed. It turns out from our measurements that the flag is seldom found clear, and thus the cost is low.

Space and time may be conserved by not logging the identity of individual slots within a chunk. It suffices simply to log which non-null references the chunk contains. This technique ameliorates the cost of spurious logging. Working on an object basis lets the mutator efficiently record precisely object slots: given a base pointer to the object we can produce a per-type slot-storing code that stores any heap slots contained in the object into the history buffer of the mutator, or produce a per-type vector of slots' indices and an efficient routine that logs the slots specified by the vector.

Allocating the dirty flag with an object is a natural choice in terms of locality, i.e., we might expect that when a slot of an object is changed, then its sibling slots are likely to change as well, so the amount of unneeded information recorded is minimized. The disadvantage of working with a flag-per-object scheme is dealing with objects that are too big. Applying the scheme for them will result in a wasteful replication of probably unchanged data. This can be avoided by treating big objects differently. Special care need be taken that the methods for small and big objects coexist.

7.2.1 An optimization. Initializing updates comprise the majority of updates in functional languages and garbage collected object oriented languages [Zorn 1990; Hosking et al. 1992]. Namely, most writes to a typical object happen just after its creation. If this is the case, we can save a substantial amount of our algorithms' overhead: when the object is created, it is logged in the mutator buffer with no contained references, and it is marked dirty. All further updates to this object up until the next collection find the dirty bit set and keep the average write barrier at low cost.

7.3 Log buffers

The primary design factor in the implementation of the log buffers is how to write into them as fast as possible for a mutator. A secondary consideration is how to allow the collector to read those records that have been fully logged (i.e., both slot and value members of a logged pair) without interrupting the mutator.

In order to satisfy the primary goal we suggest the following design, which is similar to the one described in [DeTreville 1990]: a buffer will be implemented as a linked list of *buffer-chunks*. Each chunk is of size 2^k , aligned on a 2^k boundary (k is a parameter.)

Implementing the second requirement, i.e., that the collector can read asynchronously the set of completely logged pairs can be achieved efficiently in the following manner.

- before the mutator starts using a buffer-chunk it zeroes it out.
- in order to store a record in the buffer a mutator first writes the value read, then it writes the slot address.
- the collector reads the records in the mutator's buffer sequentially. It knows that it has read a record which has not been completely logged when it sees a slot with the value of **null** (The mutator never logs a slot whose address is **null**.)

Thus, the mutator can manipulate the buffer using only a single register that points to the next address to be written.

7.4 Cooperation model - the handshakes

In order to implement handshakes, the approach in which mutators are suspended one at a time for a handshake was adopted. While a mutator is suspended, the collector takes some action on its behalf and then the mutator is resumed. We use a per-mutator flag called `cantCooperate` which is turned on in sections of code during which the mutator can not cooperate (i.e., during the write barrier and the logging of newly created objects).

In order to carry out a handshake the collector suspends the mutators one at a time. If a mutator is caught in non-cooperative code then the collector resumes it immediately and proceeds to handle other mutators. The collector repeats this process until all mutators have cooperated.

It is important to limit the size of the non-cooperative code sections to a fixed and small number of instructions. This entailed reserving space in advance, in the snoop, create and update buffers prior to entering a non-cooperative section.

8. THE MEMORY MANAGER

In the design of the memory manager we tried to satisfy these requirements:

- (1) allocation should be as fast as possible and should avoid synchronization bottlenecks. i.e., the allocator should be *scalable*.
- (2) both the tracing and reference-counting asynchronous algorithms do not accommodate the relocation of objects in memory. The allocator should not suffer from fragmentation (except maybe for some pathological cases) due to this property.
- (3) in the asynchronous reference-counting algorithm, reclamation of objects occurs sporadically rather than linearly as in the sweep phase of the tracing algorithm. The memory manager should handle efficiently this sporadic reclamation of objects. Even though objects will not be freed linearly it should still try to minimize fragmentation and increase the locality of allocation requests. i.e., it is preferable that two objects which are created in a row will be located closely in memory rather than chosen randomly from the entire heap space.
- (4) the vast majority of objects which are created are smaller than 60 bytes. The memory manager should take advantage of this fact by optimizing the allocation of small object. Allocation of medium sized and large objects may be less efficient than that of their smaller counterparts.

An allocator similar to Boehm's was chosen [Boehm et al. 1991]. The allocator is divided into two levels of management: the chunk manager and the block manager. We now outline the roles of these managers.

The block manager manages big, equally sized, blocks of memory. The block size is tunable at compile time and in the measurements was tuned to the hardware page size, which is 4KB. It supports the following operations:

- allocate a range of blocks.
- free a range of blocks given the start address of the range.
- free a collection of ranges of blocks.

The block manager is serial (No concurrency is supported) and it is implemented using linked lists of equally sized regions of blocks. The block manager is utilized either directly, by the allocation code, or indirectly using the chunk manager. When a user requests an allocation bigger than half a size of a block then the number of necessary blocks is allocated directly from the block manager. Smaller allocations are satisfied by the chunk manager which chunks single blocks into equally sized chunks that are consumed by the user.

The chunk manager is highly concurrent and efficient since it uses very fine locking, mutator local allocation and it does not support coalescing or splitting: once a block is chunked into a specific size, all allocations from within it will use the same chunk size until (and if) the block is completely freed, in which case it will be returned to the block manager. Hence, allocation code need not perform costly checks due to variable sized chunks located on the same block. There is a fixed number of allocation sizes (approximately 20). The allocation sizes are chosen to balance between internal fragmentation (which calls for many different allocation sizes) and external fragmentation (which calls for a small number of allocation sizes so that blocks of one size can be used by objects of differing sizes instead of allocating separate pages for each object size).

A typical object oriented application will issue many allocation calls that will be implemented solely by the chunk manager and only relatively few calls will require allocating entire blocks from the block manager.

9. PERFORMANCE RESULTS

In this section we report measurements we ran with our algorithm. We used two standard testing suites: SPECjbb2000 and JPECjvm98. These benchmarks are described in detail in SPEC's Web site [SPEC Benchmarks 2000].

Our primary instrumentation goal was to study the memory consumption behavior of these benchmarks. To that end, we have compiled the JVM with the GC and allocator modules in instrumented mode and the rest of the JVM in production mode. That way, the runs were realistic ones, with the amount of objects allocated and running times not significantly different from an all-production JVM yet still we gained the GC instrumentation information.

In order to appreciate the "sensitivity" of each benchmark to reference-counting, i.e., the amount of garbage cycles and stuck reference-counters that the benchmark produces, we ran each benchmark only with the tracing collector and also only with the reference-counting collector, without the use of the auxiliary tracing collector. Table I shows the number of objects allocated, average object size and the average number of references in an object. Overall, the number of allocated objects when using the reference-counting collector is comparable to the number of allocated objects using the tracing collector, though almost always smaller by a maximal factor of 5%. This is consistent with the performance figures we present later.

All tests were conducted with an equal setting for the two collectors. Our main benchmark is the standard server multithreaded application: the SPECjbb2000. This benchmark was run on a four way Pentium III at 550Mhz with 2GB of physical memory and a 600MB java heap. As a sanity check, we also ran a comparison on client applications represented by the SPECjvm98 benchmark suite. These tests were run on a single Pentium III at 500Mhz with 256MB of physical memory and

Benchmark	Tracing			RC		
	No. allocated objects	Object size	No. References	No. allocated objects	Object size	No. References
jbb	26,753,615	49.9	1.6	25,113,179	52.4	1.7
compress	55,126	2,041.1	0.8	58,061	1,940.4	0.9
db	3,261,467	34.0	2.6	3,263,358	34.0	2.6
jack	6,919,637	40.3	1.7	6,917,102	40.3	1.6
javac	6,403,821	42.9	1.9	6,405,478	43.0	1.9
jess	7,994,215	46.4	3.6	7,993,946	46.4	3.6
mpegaudio	65,539	31.6	1.1	58,329	29.5	0.9

Table I. Number of allocated objects, average object size and the average number of references in an object.

Benchmark	Stuck objects	Relative percentage
jbb	141,141	0.6%
compress	2,727	4.7%
db	30,637	0.9%
jack	51,607	0.7%
javac	235,605	3.7%
jess	12,566	0.2%
mpegaudio	2,728	4.7%

Table II. Number of objects that have reached a stuck count (i.e., 3) and their percentage in the reference-counted runs.

64MB for the jvm98 client benchmarks. The reference-counted runs of the compress and javac benchmarks were not able to complete with 64MB heap without resorting to the tracing collector to reclaim cyclic structures. Therefore, the *instrumentation* results presented here refer to runs of these two benchmarks with a java heap of 200MB. The efficiency measurements were run without this space overhead.

As Table I shows, the small number of references per object (e.g., a reference or two in a typical object) supports our premises that the number of references in most objects is relatively small hence the use of a flag per object instead of a flag per slot does not involve a significant amount of extra logging.

Table II shows the number of objects that have reached a stuck count (i.e., $o.RC = 3$) in the reference-counted runs and the relative percentage of these objects in the entire object population. These numbers support our assumption that a two-bit reference-count is enough for the striking majority of objects.

In an attempt to measure the sensitivity of each benchmark to reference-counting we compared the ratio of collected to allocated objects between the tracing and reference-counting collectors. The results are summarized in Table III. Except for javac, which uses many cyclic structures, and to a lesser degree the db benchmark, the benchmarks have demonstrated a low degree of sensitivity to reference-counting. This finding supports the premise that we may use reference-counting for most garbage collection cycles and only occasionally resort to tracing.

We now turn our attention to the use of the write barrier. Table IV shows the number of reference stores that have been applied to “new” vs. “old” objects (i.e., objects that still haven’t undergone a collection cycle versus those which have

Benchmark	% Reclaimed by tracing	% Reclaimed by RC	% Not Reclaimed by RC
jbb	97.5%	96.5%	1.0%
compress	73.5%	72.1%	1.9%
db	99.6%	90.5%	9.1%
jack	99.6%	96.8%	2.8%
javac	99.6%	66.1%	33.6%
jess	99.8%	99.5%	0.3%
mpegaudio	74.2%	69.6%	6.2%

Table III. Percentage of objects reclaimed by (1) the original collector and (2) the reference-counting collector when used without the backing mark-and-sweep collector. Ratio of percentage of objects not collected by the reference-counting collector compared with the percentage of objects collected by the tracing collector.

Benchmark	No. stores to new objects	No. stores to old objects	No. object log actions	No. references logged	Create vs. log ratio
jbb	61,070,693	9,940,664	52,410	264,115	0.00209
compress	63,892	1,013	13	51	0.00022
db	31,297,167	1,827,613	36	30,696	0.00001
jack	135,013,882	160,893	824	1,546	0.00012
javac	21,774,697	267,331	189,395	535,296	0.02946
jess	26,206,218	51,889	544	27,333	0.00007
mpegaudio	5,517,487	308	12	51	0.00021

Table IV. Demographics of the write barrier: number of reference stores applied to new and old objects; number of object logging actions; total number of references that were logged and the ratio of the number of object logging actions to the number of allocations. This ratio is an upper bound to the percentage of objects which ever get logged in the write barrier.

survived at least one collection cycle), the number of object logging actions, and the ratio of logging actions to object creation actions (this is an upper bound for the percentage of objects which ever get logged). We learn from these results the following:

- Most reference stores are applied to new objects, probably because there are more of them compared to old objects and because new objects have to be initialized.
- From the reference stores which are applied to old objects only a fraction leads to logging. This means that the same old objects are accessed repeatedly. Yet we have to log the object only the first time it is accessed in a cycle.

To conclude, due to this essentially generational behavior it is indeed beneficial to mark new objects as dirty. Also, the price paid for the write barrier almost always equals the price of a memory load and register test. Due to the large amount of new objects versus old, logged, objects, the complexity of a reference counted cycle is in reality proportional to the number of objects that were allocated during the cycle. It does appear, though, that those old objects which are repeatedly changed contain much more references compared to the average. See for example the ratio between the number of logged references to the number of logged objects in the jbb, db and jess benchmarks which far exceeds the average number of references

	jbb	compress	db	jack	javac	jess	mpegaudio
No. cycles	7	2	4	9	6	10	2
GC time	46.1	0.1	5.7	10.3	9.0	17.0	0.1
Clear	12%	7%	7%	10%	11%	7%	0%
Update	36%	19%	37%	31%	43%	37%	19%
Create buff	8%	7%	13%	13%	11%	8%	7%
Reclaim	42%	21%	41%	41%	33%	43%	21%

Table V. GC time for the reference-counting collector, in seconds. “Clear” refers to procedure **Clear_Dirty_Marks**; “Update” refers to **Update_Reference_Counters**; “Create buff” refers to the pass over the create buffers, checking whether an object is garbage and adding it to the ZCT; “Reclaim” is the final pass over the ZCT, when objects are deleted recursively.

		Heap Size (MB)	
		600	1200
score in JBB's throughput units	Original	1,131.3	1,101.0
	RC	1,101.7	1,108.3
Change in JBB score		-2.6%	0.7%
Maximal response time (milliseconds)	Original	7763	16,100
	RC	115	110
Times RC is more responsive		×67.5	×146.4

Table VI. Throughput and latency of the reference-counting collector and the original collector in standard SPECjbb runs, with 600 MB and 1200 MB heaps.

per objects in these benchmarks. This suggests that we might need to explore ways to log large objects “by pieces” and not in their entirety, as is currently done.

Finally, let us look at the execution times of each of the collectors. Table V shows the number of collection cycles, total elapsed time of the collection cycles and how this time distributes between the major stages of the reference counting collector.

For benchmarks that deal with smaller amounts of larger objects, such as compress, we see that most GC time is spent in garbage collection overheads (handshakes, etc.)

We now turn to investigate the collectors’ performance results compared to the original JVM. We start with server performance and then continue with client performance.

9.1 Server performance

A standard execution of SPECjbb requires a multi-phased run with increasing number of mutators. Each phase lasts for two minutes with a ramp-up period of half a minute before each phase. Prior to the beginning of each phase a synchronous GC cycle may or may not occur, at the discretion of the tester. We decided not to perform the discretionary initial synchronous garbage collection cycle as we believe it does not possess the characteristics of real world scenarios in which the server is not given a chance for this “offline” behavior so often. The results presented here are averaged over three standard runs.

Table VI shows the two most important performance meters for the reference-counting collector compared to the original JVM: while we do pay a small price of up to 2.6% decreased throughput, we improve the maximal response time by two orders of magnitude. To illustrate, the original JVM may pause for as long as 16

Heap size (MB)	Elapsed GC time	% increase. in GC time	No. sync cycles	No. RC cycles	No. tracing cycles
600	147	227%	2	11.7	1.0
900	144	269%	2	6.3	0.0
1200	143	225%	2	5.0	0.0

Table VII. Elapsed time of garbage collection in a standard SPECjbb run with the reference-counting collector; the percentage of increase in elapsed time over the original garbage collector and the types of garbage collection cycles that were performed. “sync” is a synchronous GC cycle requested explicitly by the benchmark.

Mutators	1	2	4	6	8	10	15	20
Original	637	1125	1728	963	928	903	887	847
RC	0.4%	4.0%	-5.4%	-2.0%	-1.0%	-2.2%	-0.3%	2.4%

Table VIII. Scores of the original JVM on a series of fixed number of mutators runs with 600MB heap; increase/decrease in score for the reference-counting collectors.

Mutators	1	2	4	6	8	10	15	20
Original	645	1137	1742	978	947	918	858	893
RC	-1.3%	3.2%	-3.8%	-3.3%	-3.3%	-3.3%	3.2%	-4.0%

Table IX. Scores of the original JVM on a series of fixed number of mutators runs with 900MB heap; increase/decrease in score for the reference-counting collector.

seconds while we never cause a mutator to pause for more than 130 milliseconds. This problem of the original JVM becomes aggravated as the heap grows in size. As can be seen from Table VII, the reason for the performance penalty is the prolonged elapsed time of garbage collection, compared to the original JVM. This implies that by further optimizing the collector code we may obtain better scores than the original JVM while maintaining the very short response time.

Next we seek to check how our collector performs relative to the original collector as a function of the number of mutator mutators and heap size. We have performed a series of stand-alone SPECjbb runs with 1, 2, 4, 6, 10, 15 and 20 mutators; 600MB, 900MB and 1200MB heaps; the original and reference-counting collector. The results are summarized in Tables VIII through XIII, and depicted in Figure 17. From throughput perspective, our collector has compatible performance with that of the original collector. We do see a slip in performance in the range of 4 to 10 mutators and this effect worsens as the heap grows. This is probably related to two factors: inefficient reclamation, which worsens as the heap grows, and tuning of spin locks for these numbers of mutators. Examining the maximal response time we again see a remarkable behavior of our collectors where the original collector consumes longer and longer pause times as the heap grows. Table XIII might seem an exception to this rule at first glance but actually what happens is that since garbage collections with such a large heap are scarce (one or two in a run) they actually might occur when the benchmark is not measuring response time hence the original JVM manages “to get away” with its long pause times unnoticed on most cases. However, examining the pause time for 4 and 20 mutators we see that these pauses nonetheless occur.

Mutators	1	2	4	6	8	10	15	20
Original	629	1155	1683	935	908	884	882	870
RC	-2.5%	1.7%	-7.1%	-8.0%	-7.8%	-6.8%	0.2%	-0.8%

Table X. Scores of the original JVM on a series of fixed number of mutators runs with 1200MB heap; increase/decrease in score for the reference-counting and tracing collectors.

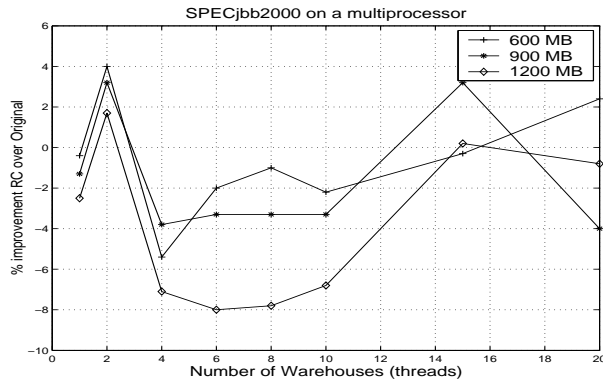


Fig. 17. Improvement in scores over the original JVM on a series of fixed number of mutators runs with 600MB, 900MB, and 1200MB heaps; increase/decrease in score for the reference-counting and tracing collectors.

Mutators	1	2	4	6	8	10	15	20
Original	7.43	8.04	8.47	6.92	7.86	7.54	6.59	6.00
RC	0.02	0.02	0.05	0.08	0.11	0.15	0.25	0.33

Table XI. Maximal response time, in seconds, of the original JVM and the reference counting collector in a series of fixed number of mutators runs with 600MB heap.

Mutators	1	2	4	6	8	10	15	20
Original	0.02	11.17	12.07	10.70	10.53	10.30	9.82	9.23
RC	0.02	0.02	0.05	0.08	0.11	0.14	0.23	0.34

Table XII. Maximal response time, in seconds, of the original JVM and the reference counting collector in a series of fixed number of mutators runs with 900MB heap.

Mutators	1	2	4	6	8	10	15	20
Original	0.02	0.02	14.67	0.05	0.08	0.01	0.18	13.03
RC	0.02	0.02	0.05	0.07	0.11	0.15	0.22	0.32

Table XIII. Maximal response time, in seconds, of the original JVM and the reference counting collector in a series of fixed number of mutators runs with 1200MB heap.

Mutators	1	2	4	6	8	10	15	20
Original	24	39	70	100	139	160	236	312
RC	27	44	77	108	170	171	251	329

Table XIV. MB allocated for heap objects and their headers (does not include space allocated for auxiliary data structures) at the end of a SPECjbb run with a fixed number of mutators and a heap of 600 MB.

Mutators	Time to completion (seconds)		% Improvement
	Original	RC	
1	93	88.6	4.9%
2	71.9	68.5	5.0%
3	56.3	52.5	7.2%
4	57.2	54.2	5.6%
8	58.2	52.3	11.4%
12	58	57.9	0.2%
16	59	59.1	-0.1%

Table XV. Time to completion, in seconds, of the MTRT benchmark, with varying number of mutators.

We now examine our memory consumption behavior. Given that we have added an extra pointer to each object (the log pointer) we would expect to see some increase in the memory consumption, relative to the average object size in each benchmark. Furthermore, since we do not compact the heap we are more vulnerable to internal fragmentation compared to the original JVM. When our collector is asked to report the amount of free memory it sums up (non-atomically) the amount of storage available in the block manager and in partial blocks. It ignores owned blocks so actually the amount of free memory is larger than reported. Given this metric, the results of used memory as reported by SPECjbb (for the 600MB test series) are summarized in Table XIV. Except for an unexplained (yet reproducible) bump in the memory consumption for 8 mutators with the reference-counting collector⁵ we consume no more than 8% more memory compared to the original JVM. This can be further improved once we eliminate completely the handle-to-object pointer in each object, which is not required by our collectors.

The second benchmark that we have used is MTRT (multi-threaded ray tracer), a member of SPECjvm98 which can be used with a varying number of mutators. We have ran this benchmark with the default heap size—64MB. This benchmark does not measure response time, only elapsed running time, which corresponds to the JVM's throughput. As can be seen from Table XV the on-the-fly collectors has outperformed the original JVM with an improvement of up to 12.6% in the total running time. The improvement is depicted in Figure 18. The reason to the decrease in improvement with a large number of mutators is that the collector is not able to finish the collection before the mutators get stuck waiting for allocation. Therefore, it becomes a stop-the-world collector. An interesting avenue for future

⁵This bump cannot be explained by reference-counting issues since the amount of consumed memory is calculated only after the benchmark requests a synchronous garbage collection cycle, which is always implemented by our collectors using a tracing cycle.

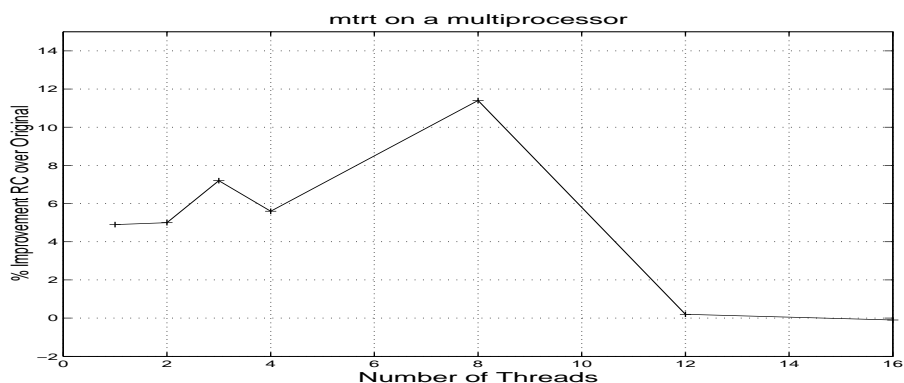


Fig. 18. Percent decrease in time to completion, of the MTRT benchmark, with varying number of mutators. Improvement in scores over the original JVM on a series of fixed number of mutators runs with 600MB, 900MB, and 1200MB heaps; increase/decrease in score for the reference-counting and tracing collectors.

	Heap (MB)	Time (sec)
Original	25	120
RC	20	85

Table XVI. Minimal heap size required to complete successfully a four thread mtrt run and the time to completion with that heap size.

research is to design a collector that is both parallel and incremental as has been done in [Ossia et al. 2002] for a tracing collector.

The ordinary measure of heap consumption—probing the free space left at the run does not capture transient effects and the ability to handle stressful situations. Table XVI shows the minimal heap size (in 1MB granularity) required to complete the mtrt benchmark successfully and the corresponding time to completion. The concurrent collectors require about 20% less the memory to complete successfully and arrive at completion at about 70% the time. This is probably a defect of the original JVM as it should actually require no more memory than our collectors and since in this stressful situation we resort to synchronous GC there should be no gain from concurrent collection as well.

9.2 Client performance

While we have targeted our collectors for multi-processor environments we still wanted to verify that they are competent in a single-processor setting. To that end we have used the SPECjvm98 benchmark suite. We used the suite using the test harness, performing standard⁶ automated runs of all the benchmarks in the suite. In a standard automated run, each benchmark is ran twice and all benchmarks are ran on the same JVM one after the other. Table XVII shows the elapsed time of the entire automated run and the time for each double run of each benchmark. We see

⁶The standard run requires running the harness through a Web server while we performed the tests directly off the disk. Aside from that, the executions were standard.

Benchmark	Original	RC
Total	2582.2	2676.0
compress	720.8	723.3
db	374.0	383.7
jack	264.6	299.7
javac	225.0	235.2
jess	181.7	209.7
mpegaudio	607.1	610.6

Table XVII. Elapsed time for the execution of the entire SPECjvm98 suite and intermediate execution time of a double-run for each of the suite’s members.

that the reference-counting collector was only 3.6% slower than the original JVM. Given that we pay the overheads of concurrent run while we’re not benefiting from the availability of multiple processors these are remarkably good results.

10. SNAPSHOT ALGORITHM CORRECTNESS PROOFS

This appendix contains safety and progress proofs for the snapshot algorithm. We start with a few necessary definitions and assumptions.

10.1 Definitions and assumptions

Global roots. For convenience, we assume in the exposition of the algorithms and their proofs that there are no global roots. Instead, we model global roots as members of a distinguished heap object which is reachable from the local state of every mutator.

Global state and time. We assume sequential consistency. thus, all shared-memory operations requested by all threads (i.e., both mutators and the collector) during a run are interleaved into a single linear order by the shared-memory system. In Section 6 we show how to adapt the algorithms to systems with weaker memory models. Assuming sequential consistency allows us to conveniently define global state and time as follows:

DEFINITION OF TIME. *For a given execution, we say that a shared-memory operation occurs at time t if it is operation number t in the linear sequence of shared memory operations corresponding to the execution.*

DEFINITION OF STATE. *For any expression E which depends only on the values of shared-memory locations and for any time point t in the execution, we denote by $E@t$ the value of entity E at time t . i.e., $E@t$ is the value of E just prior to the execution of instruction number t .*

Reachability. A mutator can access an object only if it has a local reference to it. A mutator can obtain a reference to an object only by one of two methods: (1) by reading the contents of a slot of an object to which it already has a local reference. (2) by allocating a new object. This pattern of access calls for the following standard definition of *reachability*:

DEFINITION OF REACHABILITY. *We say that an object o is*

—**directly reachable from mutator T_i at time t** *if T_i has a local-reference to o at t .*

- reachable from mutator T_i at time t** if it is directly reachable from mutator T_i at t or there exists a reference to o in object y at time t and y is reachable from mutator T_i at time t .
- reachable at time t** if there exists a mutator T_i such that o is reachable from T_i at time t .
- unreachable, or garbage, at time t** if it is not reachable at time t .

Reference-counters. Garbage collection by reference-counting is based upon counting the number of references referring to each object at a given time. We formally define the reference-count of an object as follows:

DEFINITION OF HEAP REFERENCE-COUNT. *The Heap Reference-Count of an object o at time t , denoted by $RC(o)@t$, is the number of heap slots referring to o at time t .*

We usually abbreviate and refer to an object Heap Reference-Count as its *Reference-Count*⁷. In any reference-counting system there is a field associated with each object that is used to record the number of references to the object. For an object o this field is denoted by *o.rc*. The field is invisible to the user program; it is only accessible to the memory management subsystem.

10.2 Safety

In this section, we will prove that the algorithm recycles an object only if it is garbage at the time it is recycled. Actually, an object is recycled only if it garbage at the time the conceptual snapshot is taken. Let us first define precisely this moment at which the conceptual snapshot R_k is taken:

DEFINITION 10.5. *Let HS_k be the earliest time at which all dirty marks have been cleared during the execution of procedure **Read-Current-State** in collection cycle number k . Let R_k denote the state of the heap at that moment.*

We assume that at system initialization, before any mutator has taken any step, there occur two initial garbage collection cycles. As can easily be seen, these cycles leave all data structures that are carried across cycles (e.g., reference-counters, ZCT) untouched. Since for the basis of induction proofs we use the first “real” cycle there is no loss of generality in our assumption. So, HS_{-1} and HS_0 happen at system initialization. HS_1 is the halt time corresponding to the first true collection cycle.

Ultimately, in terms of safety, we would like to prove the following:

THEOREM SAFETY. *An object is recycled during cycle k only if it is unreachable at HS_k .*

10.2.1 Road map for the proof. Here is a short and informal description of the assertions in the proof:

- SafetyTheorem_k*: An object is collected during cycle k only if it is garbage at HS_k .

⁷An object Reference-Count is sometimes defined as the number of references (including local references) to an object. We do not include local roots in the count. This definition is the same as presented in the context of *Deferred Reference-Counting*, see [Deutsch and Bobrow 1976].

- $L10.1_k$: Write barrier accurately records information.
- $L10.2_k$: If a slot is modified between HS_{k-1} and HS_k then only and exactly the value it assumed at HS_{k-1} is recorded. No information is recorded for slots which are not modified.
- $L10.3_k$: The collector can distinguish, during cycle k , whether it is reading a slot's value which was current at HS_k , or, that the slot has been overwritten since.
- $L10.4_k$: The collector finds out, eventually, in procedure **Fix-Undetermined-Slots**, what the values of undetermined slots are.
- $L10.5_k$: Just before the invocation of **Reclaim-Garbage** during cycle k , the rc field of each object equals the heap reference-count of the object at HS_k .

Most lemmas are interdependent meaning, for example, that we prove lemma X correct at cycle k provided lemma Y is correct at cycle $k - 1$. In order to make clear the relation between the claims and to demonstrate that there is no circular logic in the proof we provide herein a complete description of the interdependencies among the claims. We denote by Li_k the assertion of lemma i for cycle k .

These are the dependencies between the claims:

- the basis for each claim, i.e. its correctness for cycle one is proven independently for each claim.
- $L10.2_k \Leftarrow \bigwedge_{j < k} SafetyTheorem_j$
- $L10.5_k \Leftarrow \bigwedge_{j < k} (L10.5_j \wedge SafetyTheorem_j) \wedge L10.2_k \wedge L10.4_k$
- $SafetyTheorem_k \Leftarrow L10.5_k$

10.2.2 Update protocol properties. Consider any slot s which is modified between HS_{k-1} and HS_k . The snapshot algorithm requires us to adjust rc fields due to s by decrementing the rc field of $s@HS_{k-1}$ and incrementing the rc field of $s@HS_k$. The first part of the requirement, decrementing $s@HS_{k-1}$, is implemented by letting the mutators record the identity of $s@HS_{k-1}$ into their buffers. Thus, we would like to prove for any such modified slot s that only and exactly $s@HS_{k-1}$ is associated with s by the mutators.

If s is not modified between the current and previous cycles, then we want to show that no record of s is kept.

The lemmas in this section prove that the algorithm possesses these properties.

LEMMA 10.1. *Let s be a slot and let t be a time point satisfying*

- (1) $HS_{k-1} \leq t < HS_k$, and
- (2) $Dirty(s)@t = \mathbf{false}$, and
- (3) No update of s is occurring at t .

Let $UPD(s)$ be the set of all update operations applied to s which are scheduled between t and HS_k . Let $ASSOC(s)$ be the set of values written to the buffers for s by the operations in $UPD(s)$.

It holds that:

- (1) $UPD(s) = \emptyset \implies ASSOC(s) = \emptyset$
- (2) $UPS(s) \neq \emptyset \implies ASSOC(s) = \{s@t\}$

PROOF. The first claim is quite trivial since a value is associated with s only as part of an update. Since no update is scheduled, no value is associated.

Suppose that s is indeed modified between t and HS_k . Consider the set of mutators, denoted P , that apply the subset of operations of $UPD(s)$ which read the value of $Dirty(s)$ as **false** in line (2) of procedure **Update**, while updating s . P is not empty since some mutator modifies s ($UPD(s)$ is non-empty) and the dirty flag is off at t .

Consider a mutator $T_i \in P$. We want to show that when T_i executed line (1) of procedure **Update** it read the value of s at t . Suppose that it did not. Let τ be the time at which mutator T_i executed line (1). Then some mutator T_j must have executed a store to s after, or at, t and before τ . Since there were no updates occurring at t and since the store is the last instruction of an update operation we conclude that the entire update operation by T_j has started after, or at, t and ended before τ . Just before T_j executed the store in line (6) the value of $Dirty(s)$ must have been **true** either by line (5) or by virtue of another mutator (the collector resets the flag only during the next cycle) so T_i should have read a value of **true** from $Dirty(s)$, in line (2), which was not the case. A contradiction. We conclude that T_i must have associated $s@t$ with s . So we have

$$\{s@t\} \subseteq ASSOC(s)$$

According to the code, any mutator $T_i \notin P$ would not associate any value with s thus

$$ASSOC(s) = \{s@t\}$$

. Additionally, values are copied from the buffers without duplications. Thus, the same value may not appear twice. \square

For a given history buffer H (be it collector or mutator maintained set) and a slot s we define the set of values that H associates with s , denoted by $VAL(H; s)$, as:

$$VAL(H; s) \stackrel{\text{def}}{=} \{v \mid \langle s, v \rangle \in H\}$$

For brevity we write $s \in H$ meaning $\exists v : \langle s, v \rangle \in H$. The next lemma summarizes and proves the desired properties of the write-barrier employed by the algorithm. We need some definitions first:

- We say that an object o is *allocated for cycle k* . If some mutator has allocated o between HS_m and HS_{m+1} , where $m < k$. And there has not been a cycle l , where $m \leq l < k$ during which o was reclaimed.⁸
- o is *allocated new for cycle k* if $m = k - 1$ in the above definition.
- If $m < k - 1$, we say that o is *allocated old for cycle k* .
- We say that a slot is *allocated (new/old) for cycle k* if its containing object is allocated (new/old) for cycle k .
- We abbreviate and say that a slot or an object are *new (old) to a cycle* meaning that the slot or the object are allocated new (old) for that particular cycle.

⁸Even though l never equals m in the algorithm, we don't want to assume that *a priori*.

LEMMA 10.2. *Let s be an allocated slot for cycle k . Then:*

(1) *if s is new to cycle k and is modified between HS_{k-1} and HS_k then*

$$VAL(Hist_k; s) = \{\mathbf{null}\}$$

(2) *if s is old to cycle k and is modified between HS_{k-1} and HS_k then*

$$VAL(Hist_k; s) = \{s@HS_{k-1}\}$$

(3) *otherwise (s is not modified between HS_{k-1} and HS_k),*

$$VAL(Hist_k; s) = \emptyset$$

PROOF. The lemma vacuously holds for $k = 1$ since there are no slots which are modified during the interval HS_0 to HS_1 .

We now show that the lemma holds for cycle $k > 1$ provided that the safety theorem hold for previous cycles.

Suppose s is new to cycle k . Let τ be the time at which the object o containing s was allocated. Let $j < k$ be the cycle during which the object x that most recently contained s was reclaimed, or 0 if no such cycle exists. Applying the safety theorem to cycle j we know x was unreachable at HS_j . Thus, no mutator could have accessed s from HS_j until τ . If $j > 0$ then when x was recycled a **null** value was assigned to s , in line (4) of procedure **Collect**; otherwise, s contained **null** since system startup. So at any rate s assumed the value of **null** at HS_j . Also, as all dirty flags are cleared while the mutators are halted, we have $Dirty(s)@HS_j = \mathbf{false}$. Since the values of s and $Dirty(s)$ must remain constant until time τ and since no mutator is updating s at time τ (all due to the safety theorem applied to cycle j) we can apply Lemma 10.1 to s and τ yielding that either claim (1) or (3) holds, depending on whether s has been modified between τ and HS_k .

If, on the other hand, s is old to cycle k then we have $Dirty(s)@HS_{k-1} = \mathbf{false}$ and no update of s is occurring at HS_{k-1} . Thus, we can apply Lemma 10.1 to s and time HS_{k-1} yielding that either claim (2) or (3) hold, depending on whether s has been modified prior to HS_k . \square

10.2.3 *Determined vs. undetermined slots.* We say that the collector *determines* the value of a slot s if during the **Update-Reference-Counters** procedure it reads the value v from s (in line (3)) and then sees $Dirty(s) = \mathbf{false}$ (in line (4)). Such a slot is *determined*, as opposed to *undetermined* slots which are taken care of by the collector in procedures **Read-Buffers** and **Fix-Undetermined-Slots**. The following lemma tells us that if the collector determines the contents of a slot then it has indeed read its contents as they were at the time the recent conceptual snapshot was taken.

LEMMA DETERMINED SLOTS. *If the collector determines s to contain v during cycle k then $v = s@HS_k$.*

The proof is similar to that of Lemmas 10.1 and 10.2.

What happens when the collector does not succeed determining a slot? A slot is undetermined if the collector senses that its flag is raised during **Update-Reference-Counters**. The only reason for the flag to be raised is that some mutator, say T_i , has applied line (5) of procedure **Update** to the flag (i.e., raised

it.) In this case, T_i has executed the preceding lines of (3) and (4) of the same invocation after HS_k . i.e., T_i has stored the pair $\langle s, s@HS_k \rangle$ into its buffer and incremented $CurrPos_i$ prior to raising the flag. Thus, when the collector would process $Buffer_i$ during **Read-Buffers** it will see the logged pair $\langle s, s@HS_k \rangle$ in T_i 's buffer ($s@HS_k$ is associated with s according to Lemma 10.2.) and thus the pair will be added to the set $Peek_k$.

We conclude the following:

LEMMA UNDETERMINED SLOTS. *If the collector does not determine a slot s in cycle k then*

$$VAL(Peek_k; s) = s@HS_k$$

10.2.4 *Linking rc field with reference-count.* In this section we show that the rc fields that the algorithm computes equal, eventually, the heap reference-counts at the time the conceptual snapshot is taken. We need some definitions first.

DEFINITION 10.6. *Let END_k denote the time at which cycle k has ended. That is, END_k is the earliest time at which all instructions of cycle k have already been scheduled.*

DEFINITION 10.7. *Let $COLLECT_k$ be the time at which the invocation of **Fix-Undetermined-Slots**, during cycle k , is complete. The collector starts executing **Reclaim-Garbage** after, or at, $COLLECT_k$.*

The following lemma asserts that the value of the rc field of each object, after the collector has finished adjusting rc fields due to all logged modifications, i.e., when procedure **Reclaim-Garbage** starts its operation, equals the object's heap reference-count at time HS_k .

LEMMA MEANING OF THE rc FIELD. *$o.rc@COLLECT_k = RC(o)@HS_k$ for any object o which is allocated at HS_k .*

PROOF. The claim holds for $k = 0$ since there are no objects which are allocated at HS_0 . For $k > 0$, we prove that the lemma holds for cycle k provided this lemma and the safety theorem both hold for previous cycles.

It is enough to show that the algorithm adjusts rc fields due to each slot s correctly. If s does not change after HS_{k-1} and before HS_k then, by Lemma 10.2, s will not be logged and there will be no modifications to any rc fields due to s .

Let's consider the cases in which s does change. We have to show that the rc field of the object that s was referring to at HS_{k-1} is decremented. Likewise, we have to show that the value of the object that s was referring to at HS_k is incremented. s is in exactly one of these states at HS_k : allocated old, allocated new, non-allocated.

Decrementing old slots: If s is old for cycle k then s is changed by mutators, and not by the collector (by deleting it.) Due to Lemma 10.2 $Hist_k$ will contain the pair $\langle s, s@HS_{k-1} \rangle$. $Hist_k$ will not contain elements associating s with a value other than $s@HS_{k-1}$. During the operation of **Update-Reference-Counters**, when the pair $\langle s, s@HS_{k-1} \rangle$ is considered, the rc field of $s@HS_{k-1}$ is decremented, as desired.

Decrementing new slots: Let s be a new slot for cycle k . According to Lemma 10.2 either **null**, or no value at all, are associated with s . Thus, there are

no decrements that occur due to s during cycle k . Let us explain why this is the desired behavior.

If s is new for cycle k then either s becomes allocated for the first time, or it was part of an object o which was recycled during cycle j , where $j < k$.

In the former case, we know that s was initialized to **null** and its dirty flag was off at system startup. Also, no mutator could have accessed s at HS_{k-1} , since it was not a part of a reachable object (or any object) at that time. Thus, $s@HS_{k-1} = \mathbf{null}$ and therefore no rc field should be decremented due to s during cycle k .

In the latter case, according to the safety theorem applied to cycle j , o is not reachable at HS_j . Thus, the collector has exclusive access to s , during cycle j . It follows that the collector may decrement the rc field of the object referenced by s and clear s without being interfered by mutators' actions, all part of the operation of **Collect** during cycle j . If $j < k-1$ then $s@HS_{k-1} = \mathbf{null}$, thus there is no "old" value to decrement.

Otherwise, $j = k-1$. In this case, the collector decrements the rc field of $s@HS_{k-1}$ during cycle $k-1$, when it reclaims o . An object is reclaimed only if its rc field drops to zero. **Reclaim-Garbage** and **Collect** can only reduce the value of an rc field. Thus, there is a single point during the operation of **Reclaim-Garbage** at which $o.rc = 0$. Therefore o is reclaimed exactly once and likewise the rc field of $s@HS_{k-1}$ is decremented exactly once.

Decrementing and incrementing non-allocated slots: If s is not allocated at HS_k then the same argument that was applied to new slots is used to show that the value of $s@HS_{k-1}$ is taken care of. Again, due to the safety theorem applied to the cycle at which the object containing s was recycled we have $s@HS_k = \mathbf{null}$ so there is no need to increment any rc field due to s . Indeed, since s is not allocated at HS_k and it is unreachable at HS_{k-1} no record of it would appear in $Hist_k$ and no rc field will be manipulated due to it in cycle k .

Incrementing old and new slots: it remains to show that the rc field of $s@HS_k$ is incremented exactly once due to s , when s is allocated at HS_k . We have two cases: either s is determined, or it is undetermined. If s is determined, then due to Lemma 10.3 we have that the collector increments the rc value of $s@HS_k$. Otherwise, by Lemma 10.4, $VAL(Peek_k; s) = \{s@HS_k\}$. Thus, during the **Fix-Undetermined-Slots** procedure the collector will find the value of $s@HS_k$ associated with s . It will increment the rc field of that object exactly once, by the code.

All rc adjustments are finished by the time **Fix-Undetermined-Slots** terminates, so the claim holds at $COLLECT_k$. \square

10.2.5 *Conclusion of safety proof.* We are now ready to prove the safety theorem which claims that an object is collected at cycle k only if it is unreachable at time HS_k .

PROOF. (of safety theorem.) The claim trivially holds for cycle zero since ZCT_0 is an empty set and thus no object is recycled during the initial cycle.

Consider cycle $k > 0$. We prove that the theorem holds for cycle k if Lemma 10.5 holds for cycle k .

Let $\{o_1, \dots, o_n\}$ be the sequence of objects for which **Collect** is invoked, where the sequence is chronologically ordered. We show by induction on i , that o_i is unreachable at HS_k . For the basis, consider o_1 . As it is the first object to be collected, there is no clearing of slots (carried out in line (4) of procedure **Collect**) taking place prior to its reclamation, thus $o_1.rc@COLLECT_k = 0$. This implies, according to Lemma 10.5 applied to cycle k , that $RC(o)@HS_k = 0$. Additionally, by the code, o_1 is collected only if $o_1.rc = 0 \wedge o_1 \notin Locals_k$ so we conclude that in addition of not being referenced by any heap slot at HS_k , o_1 is also not referenced by any local reference at that particular moment, or it would have been marked *local*. Thus, o_1 is unreachable at HS_k .

For the inductive step, consider o_i which has $c \stackrel{\text{def}}{=} o_i.rc@COLLECT_k = RC(o)@HS_k$ (the last equality is again by Lemma 10.5). If $c = 0$ then the same arguments that were employed for o_1 are repeated in order to demonstrate that o_i is garbage at HS_k .

Otherwise, we have $c > 0$. Since o_i is recycled, it must satisfy at some point during **Reclaim-Garbage** or **Collect** $o_i.rc = 0 \wedge o_i \notin Locals_k$. Thus, the value of $o_i.rc$ is decremented c times during the operation of **Reclaim-Garbage**. Since decrements are only applied to objects which are referenced by objects that are collected and since those objects are collected prior to o_i we have by the inductive hypothesis that all c references to o_i were from objects that were unreachable at HS_k . Thus, at HS_k , o_i is referenced only by unreachable objects, and it is not referred to by any local mutator state or global reference. We conclude that o_i is unreachable at HS_k . \square

10.3 Progress

In this section we show the capabilities of the algorithm in collecting garbage objects. The algorithm, in that respect, has the same limitations as the traditional single-threaded reference-counting algorithms [McBeth 1963].

The best that we can hope to achieve with reference-counting, without employing special techniques for detecting cycles of garbage, such as those surveyed in [Lins and Vasques 1991], is to detect any object that its reference-count drops to zero, in order that it would be considered for reclamation based on the existence of local references to it. The following lemma tells us that this feature is achieved by the ZCT data-structure.

LEMMA ZCT PROPERTY. *If o is allocated at HS_k and $RC(o)@HS_k = 0$ then $o \in ZCT_k$.*

PROOF. The proof is by induction on k . There are three cases to consider:

- (1) o is new to cycle k . In this case, a mutator created o between HS_{k-1} and HS_k . When it created o it added it to its *New* set, which becomes part of ZCT_k .
- (2) o is old to cycle k and it had a positive *rc* field at END_{k-1} . Since we have $0 = RC(o)@HS_k = o.rc@COLLECT_k$ (by Lemma 10.5), the value of $o.rc$ must have reached zero due to the decrements applied by procedure **Update-Reference-Counters** of cycle k . At that point o was added to ZCT_k (see lines (8-10) of that procedure.)

(3) o is old at HS_{k-1} and it had zero rc field at END_{k-1} . This case splits into two sub-cases:

- (a) if $o.rc@COLLECT_{k-1} = 0$ then $RC(o)@HS_{k-1} = 0$ by Lemma 10.5. Using the inductive assumption we know that $o \in ZCT_{k-1}$. Since o was not recycled we must have $o \in Locals_{k-1}$. By the code, when o is considered during **Reclaim-Garbage** it satisfies

$$o.rc = 0 \wedge o \in Locals_{k-1}$$

by the code (lines (5-7)), o is added to ZCT_k in this case.

- (b) Otherwise, $o.rc@COLLECT_{k-1} > 0 \wedge o.rc@END_{k-1} = 0$. This implies that $o.rc$ had reached zero by the decrements applied by one of the invocations of procedure **Collect**. By the code (lines (5-9)), when an object reference-count reaches zero but it is not reclaimed, it is moved to the ZCT of the next cycle.

□

Ideally, we would like the algorithm to collect at cycle k any object which is garbage at HS_k . However, this algorithm has the standard weakness of reference-counting, with respect to cyclic structures, and thus only the following progress theorem can be guaranteed:

THEOREM PROGRESS. *If at HS_k object o is unreachable and additionally o is not reachable from any cycle of objects, then o is collected in cycle k .*

The theorem follows from Lemma 10.6 and the fact that we use standard recursive-freeing.

11. SLIDING VIEW ALGORITHM SAFETY PROOF

In this appendix we prove that the sliding view algorithm is safe.

11.1 Definitions

First we need to extend the definitions a bit in order to accommodate the looser timing of the sliding views algorithm.

Let us define the time instances at which a mutator T_i is suspended during the four handshakes of each cycle: $HS1_k(i)$, $HS2_k(i)$, $HS3_k(i)$ and $HS4_k(i)$ denote the time instances at which mutator T_i is suspended during the first, second, third and fourth handshakes of cycle k , respectively. Next, we define the “global” time markers at which each handshake starts (by stopping the first mutator) and ends (by resuming the last stopped mutator):

$$\begin{aligned} HS1_k &\stackrel{\text{def}}{=} \min_{T_i} HS1_k(i) \\ HS1END_k &\stackrel{\text{def}}{=} \max_{T_i} HS1_k(i) \\ HS2_k &\stackrel{\text{def}}{=} \min_{T_i} HS2_k(i) \\ HS2END_k &\stackrel{\text{def}}{=} \max_{T_i} HS2_k(i) \\ HS3_k &\stackrel{\text{def}}{=} \min_{T_i} HS3_k(i) \\ HS3END_k &\stackrel{\text{def}}{=} \max_{T_i} HS3_k(i) \end{aligned}$$

$$HS4_k \stackrel{\text{def}}{=} \min_{T_i} HS4_k(i)$$

$$HS4END_k \stackrel{\text{def}}{=} \max_{T_i} HS4_k(i)$$

Additionally we define $COLLECT_k$ to be the time at which procedure **Reclaim-Garbage** starts its operation. The notions of “being allocated” of the snapshot algorithm’s proof must also be modified due to the lack of the hard handshake. This is done in the following definitions:

- We say that an object o is *allocated for cycle k* if some mutator T_i allocated o after $HS1_m(i)$ but before $HS1_{m+1}(i)$, where $m < k$, and there had not been a cycle l , where $m \leq l < k$, such that o was reclaimed on cycle l .
- o is *allocated new for cycle k* if $m = k - 1$ in the above definition.
- If $m < k - 1$, o is *allocated old for cycle k* .
- We abbreviate and say that o is *new (old) to cycle k* if it is allocated new (old) for cycle k .
- Any of the above definitions applies to slots by letting the definition hold for the object containing the slot.

11.2 The sliding view associated with a cycle

In this section we define a per-cycle sliding view . Please refer to Section 5.1 for the definition of scans and views. In short, $\sigma(s)$ is the time at which the slot s is scanned, and $V_\sigma(s)$ is the value probed during the scan. The per-cycle sliding view is later shown to be computed implicitly by the collector and mutators (bearing similarity to the conceptual snapshot taken at HS_k by the snapshot algorithm which is never explicitly computed.)

Let us define the scan σ_k that we associate with cycle k . We abbreviate V_{σ_k} to V_k . Consider any memory word s . If s does not appear in the buffers that are read during the first handshake of collection k , then s does not cause any RC updates, and we may choose to think of it as being read by the sliding view scan at $HS1_k$. If s does appear in $Hist_k$, we split the discussion into two cases. If s is logged by some mutator i between $HS1(i)$ and $HS3(i)$ by some mutator T_i , then its value for $Hist_{k+1}$ will be set by the procedure **Get-Local-States**. In this case, the scan time is defined according to the consolidated value for s . Otherwise, s is not logged by any mutator between the first and third handshakes and then we choose the scan time to be $HS2END_k$. Formally, the scan is determined by the following set of rules.

- Rule 1:* if $s \notin Hist_k$ then we set $\sigma_k(s) = HS1_k$.
- if $s \in Hist_k$ then:
 - Rule 2:* if s is logged by some T_i between $HS1_k(i)$ and $HS3_k(i)$ then let v be the consolidated value chosen for s . Let τ be the time a particular mutator T_j loaded v before logging the pair $\langle s, v \rangle$. Set $\sigma_k(s) \stackrel{\text{def}}{=} \tau$.
 - Rule 3:* otherwise, no mutator T_i logs s between $HS1_k(i)$ and $HS3_k(i)$, but s is logged by some mutator T_j prior to $HS1_k(j)$. On such an event set $\sigma_k(s) \stackrel{\text{def}}{=} HS2END_k$.

We denote by $R1_k$ the set of all slots whose definition of σ_k is derived by rule (1). Similarly we define the sets $R2_k$ and $R3_k$.

The next lemma characterizes the span of σ_k .

LEMMA 11.1. $Start(\sigma_k) \geq HS1_k \wedge End(\sigma_k) \leq HS3END_k$

PROOF. Let s be a memory word. Certainly if $s \in R1_k \cup R3_k$ then $\sigma_k(s)$ lies within the specified time limits. Otherwise, s is defined according to rule (2). Since $s \in Hist_k$ it may be read and logged again only after its dirty bit is cleared, i.e., after $HS1_k$. Thus, its value is read and logged by the mutator T_j of Rule 2 after the first handshake ends. By definition of Rule 2, s has been logged prior to the third handshake. Although the consolidation process occur only after the fourth handshake, we claim that τ must be earlier than $HS3END_k$. By the definition of Rule 2, some mutator must log s prior to responding to the third handshake. If this logging is done during clearing (i.e., between the first and second handshake) then the flag will be reinforced before the third handshake. Otherwise, the flag must remain on until the clearing of the next cycle. So at any rate, the flag is on at $HS3END_k$. Thus no mutator could load a value from s after $HS3END_k$ and then log it since it is bound to sense that the dirty flag of s is on. \square

The following lemma links the asynchronous reference-count associated with the cycle's scan with the instantaneous reference-count at $HS1_k$.

LEMMA 11.2. *Let V_σ be a sliding view and let o be an object. If for any slot s , no reference to o is stored into s at, or after, $\sigma(s)$ and before $End(\sigma)$ then $RC(o)@End(\sigma) \leq ARC(V_\sigma; o)$. Furthermore, the set of slots that refer to o at $End(\sigma)$ is a subset of those that refer to it in V_σ*

The proof is a simple extension of the proof of Proposition 5.4 and is omitted.

LEMMA 11.3. *Any object o which is not marked local (i.e., $o \notin Locals_k$) at $COLLECT_k$ satisfies*

$$RC(o)@HS4_k \leq ARC(V_{\sigma_k}; o)$$

Moreover, the set of references that refer to o at $HS4_k$ is a subset of those that refer to it in V_{σ_k} .

PROOF. According to Lemma 11.2 it suffices to show that if a reference to o is stored to a slot s at, or after $\sigma_k(s)$ and before $End(\sigma_k)$, then o is marked local. By Lemma 11.1 we know that $End(\sigma_k) < HS3END_k < HS4_k$, hence it is enough to require that if a reference to o is stored to a slot s during the interval $[\sigma_k(s), HS4_k)$ then o is marked local.

Since the $Snoop_i$ flag is reset only after $HS4_k(i)$, it suffices to show that the test of $Snoop_i$ in the **Update** procedure returns **true** in the case that the store proper into s is executed after $\sigma_k(s)$ and before $HS4_k(i)$. Consider a store of o into s which is scheduled at, or after $\sigma_k(s)$ and before $HS4_k(i)$. Due to Lemma 11.1, the store is scheduled at or after $HS1_k$. At that time, for any mutator T_i , the $Snoop_i$ flag is set. Since the test of $Snoop_i$, in line (7) of procedure **Update**, is executed after the store proper, of line (6), it would return **true** and the object will be marked accordingly *local*. The fact that handshakes do not stop mutators in the middle of the write-barrier, guarantees that any store that starts before $HS4_k$ will finish

the write-barrier including proper snooping, before clearing the snoop flag during handshake 4. \square

We now turn into validating the dirty bit behavior.

LEMMA 11.4. *For any mutator T_i and any collection k the following holds.*

- (1) *if mutator T_i logs s between responding to the first and third handshakes then $Dirty(s)@HS3_k(i) = \mathbf{true}$.*
- (2) *if mutator T_i logs s between responding to the first and fourth handshakes then $Dirty(s)@HS4_k(i) = \mathbf{true}$.*

PROOF. Part (1): A dirty flag may be cleared after T_i has raised it only if the collector has reset the dirty flag in procedure **Clear-Dirty-Marks**. If that is the case, then the collector has reset the flag after the mutator T_i has completed logging the slot. Hence, in procedure **Reinforce-Dirty-Mark**, the collector will see the slot in T_i 's buffer and would reinforce it. This happens before $HS2_k$. Part (2) follows in the same manner noting that the dirty flag is only cleared between handshakes 1 and 2. \square

11.3 Inductive safety arguments

Now that a sliding view with each cycle has been defined and some properties about it have been proven, we turn to proving the safety theorem. The proof of the safety theorem is by induction on the collection cycle number.

Compensating for the lack of the hard handshake of the snapshot algorithm, during which all dirty marks were turned off we have procedure **Clear-Dirty-Marks** in the sliding view algorithm. The following lemma asserts that indeed each slot experiences a point in time, after the start of a cycle, at which the dirty flag is off. This is essential for the logging mechanism to operate correctly since it instructs mutators to start logging modifications from fresh, relating to the new cycle.

LEMMA 11.5. *Let s be a heap slot. There exists a time point, denoted $t_k(s)$ at which the dirty flag for s is off. Specifically:*

- if $s \in R1_k$ then $t_k(s) \stackrel{\text{def}}{=} \sigma_k(s) \stackrel{\text{def}}{=} HS1_k$.
- if $s \in R2_k$ then $t_k(s)$ exists and it satisfies $HS1END_k < t_k(s) < HS2_k$.
- if $s \in R3_k$ then $t_k(s) \stackrel{\text{def}}{=} HS2END_k$. There are no ongoing updates of s at $t_k(s)$.

PROOF. The proof is by induction on the cycle number, k . For $k = 0$ the claim holds since all slots are cleared at $HS1_0$ and all slots are members of $R1_0$. For $k > 0$ we prove the claim correct provided it holds for the previous cycle and Theorem 11.1 holds for all previous cycles. We divide to cases:

- if $s \in R1_k$ then either $s \in R1_{k-1}$ or $s \in R3_{k-1}$. $s \in R2_{k-1}$ is impossible because it implies that $s \in Hist_k$.

If $s \in R1_{k-1}$ then by the inductive hypothesis $Dirty(s)@HS1_{k-1} = \mathbf{false}$. Had some mutator T_i turned on the flag on after $HS1_{k-1}$ and before $HS1_k(i)$ then s would have been recorded in either $Hist_{k-1}$ or $Hist_k$, neither of which is the case, so the dirty flag must be continuously off from $HS1_{k-1}$ to $HS1_k$.

Otherwise, $s \in R3_{k-1}$. Thus, according to the inductive hypothesis $Dirty(s)$ is turned off at $HS2END$. By definition of $R3_{k-1}$, no mutator logged s before responding to the third handshake of cycle $k - 1$. Thus no mutator had turned the flag on prior to responding to that handshake. Had some mutator logged s after the third handshake of cycle $k - 1$ but before the first handshake of cycle k then we would have $s \in Hist_k$, which is not the case. Again we have $Dirty(s)@HS1_k = \mathbf{false}$.

- if $s \in R2_k$ then the collector has turned off $Dirty(s)$ during the clearing stage. We define $t_k(s)$ to be the time instance just after the clearing of $Dirty(s)$ was scheduled.
- if $s \in R3_k$ then the collector has turned off $Dirty(s)$ during the clearing stage and no mutator has turned it on prior to responding to the third handshake. We conclude that the flag must have been off at the time the second handshake ended. At $HS2END_k$ only updates of mutators that have already responded to the second handshake may be ongoing. But had such an update occurred, it must have sensed that the flag is off and it would consequently log s , contradicting the definition of $R3_k$. We conclude that there are no ongoing updates at $HS2END_k$.

□

The properties of the write-barrier are now considered. The next lemma, which is the equivalent of Lemma 10.2 of the snapshot algorithm, states that any slot which is modified between scans is recorded along with its value in the previous sliding view and that no other value is associated with the slot.

LEMMA 11.6. *Let s be a slot. The following claims hold:*

- (1) *if s is old for cycle k and modified during cycle $k - 1$ then $VAL(Hist_k; s) = \{V_{k-1}(s)\}$.*
- (2) *if s is new for cycle k and modified during cycle $k - 1$ then $VAL(Hist_k; s) = \{\mathbf{null}\}$.*
- (3) *if s is old for cycle k and is not modified during cycle $k-1$ then $VAL(Hist_k; s) \subseteq \{V_{k-1}(s)\}$.*
- (4) *if s is new for cycle k and is not modified during cycle $k-1$ then $VAL(Hist_k; s) \subseteq \{\mathbf{null}\}$.*

PROOF. For garbage collection number zero the claims trivially hold since $Hist_0 = \emptyset$ and indeed no slot is modified prior to the cycle. We prove that the claim holds for cycle $k > 0$ provided it itself hold for cycle $k - 1$ and that Theorem 11.1 and Lemma 11.5 hold for earlier cycles.

The following cases are possible according to the state of s . Old s 's may be in $R1_{k_1}$, in $R2_{k_1}$, or in $R3_{k-1}$ and are treated accordingly. New slots s are partitioned according to whether they have been allocated before or this is their first allocation. We start with old s 's.

Case 1: s is old for cycle k and $s \in R1_{k-1}$. Case 1 is further divided. Suppose that $s \notin Hist_k$. In that case we have $\sigma_k(s) \stackrel{\text{def}}{=} HS1_k$ and we have to show that s is not changed between $HS1_{k-1}$ and $HS1_k$. i.e., we have to show that claim (3) is not violated in this case, because all other claims do not apply. Since $s \notin Hist_{k-1}$ we

conclude, by the inductive hypothesis, that no mutator modified s between $\sigma_{k-2}(s)$ and $HS1_{k-1}$. Additionally we know that at $HS1_{k-1}$ the dirty mark of s is off. The dirty mark must be off at $HS4END_{k-2}$ as well and no update is ongoing at the moment as that update would have rendered s part of $Hist_{k-1}$. Using the same arguments of Lemma 10.1 applied for s and $HS4END_{k-2}$ and since s is not cleared before $HS1END_k$ any update whose store proper operation is scheduled between $HS4END_{k-2}$ and $HS1_k$ would result in the association of $s@HS4END_{k-2}$ with s in either $Hist_{k-1}$, or $Hist_k$, neither of which is the case. We conclude that s is indeed not modified during cycle $k - 1$.

Now suppose $s \in Hist_k$. In that case we want to show that $VAL(Hist_k; s) = s@HS1_{k-1}$. Again, we've concluded that any mutator T_i that would log s prior to $HS1_k(i)$ would associate it with $s@HS4END_{k-2}$. Since a store to s could not have been scheduled between $s@HS4END_{k-2}$ and $HS1_{k-1}$ without logging the slot we conclude that $s@HS1_{k-1} = s@HS4END_{k-2}$, which is the desired result.

Case 2: s is old for cycle k and $s \in R2_{k-1}$. Since some mutator modified and logged s between the first and third handshakes of cycle $k - 1$ We have to show that claim (1) holds for s . Due to the reinforcement step, the dirty flag of s must be on at $HS4_{k-1}$, thus, there is no possibility that a mutator would log s after responding to the fourth handshake. As for the records kept regarding s between the first and fourth handshakes, the collector chooses a single pair, say $\langle s, v \rangle$ and moves it to $Hist_k$. By definition of σ_k we have $V_{k-1}(s) = v$.

Case 3: s is old for cycle s and $s \in R3_{k-1}$. We have noted in Lemma 11.5 that $t_{k-1}(s) = \sigma_{k-1}(s) = HS2END_{k-1}$ and no update is occurring at that moment. Suppose $s \notin Hist_k$. In that case $\sigma_k(s) = HS1_k$ and we have to show that no store is scheduled between $HS2END_{k-1}$ and $HS1_k$. But this is trivial since the probing of the dirty mark associated with such a store must start after $HS2END_{k-1}$, as no updates occur at that moment. Thus, had such an update been scheduled, it must have sensed that the flag is off and s would have become a member of $Hist_k$ a contradiction.

Suppose now that $s \in Hist_k$. We have to show that $VAL(Hist_k; s)$ equals $s@HS2END_{k-1}$. Again, since at $HS2END_{k-1}$ the dirty bit is off and no update of it is occurring. And since the dirty mark is reset only after all mutators have responded to the first handshake of cycle k , by Lemma 10.1 they are bound to associate $s@HS2END_{k-1}$ with s .

Case 4: new slots allocated for the first time. If s is allocated for the first time, then $\sigma_{k-1} \stackrel{\text{def}}{=} HS1_{k-1}$ and at that time s contained **null** and its dirty flag was initialized to **false**. These values remain in effect until s is allocated. Additionally, no update of s occurs at the moment it is allocated. Again, the claim follows using the arguments of the previous cases.

Case 5: new slots which are reallocated. We first show that $Hist_k$ cannot contain “leftovers”: i.e., logging that refer to the “previous life” of s , before it was reallocated. Suppose that s was last reclaimed during cycle m , $m < k$. If $m < k - 1$, then there will be no record of the “previous life” of s in $Hist_k$ due to the safety theorem applied to cycle m that assures us that s was unreachable from its reclamation point up to the time it was re-allocated, during cycle $k - 1$. If, on the other hand, s was reclaimed during cycle $k - 1$, then as the safety theorem tells us, no mutator T_i had access to s after $HS4_{k-1}(i)$. s could have not occurred in the

digested part of $Hist_k$ as that would have caused the deferral of the reclamation of its containing object to cycle k . So there are no leftovers in this case as well.

Applying the safety theorem to cycle m , we know that the object that contained s was garbage when it was reclaimed. Its dirty marks, the one of s included, were off. When the collector freed the object it stored **null** into s . Since the object was unreachable, s remained inaccessible up to the time it was re-allocated. Just when s was re-allocated, there was no update of it ongoing, it contained **null**, and the dirty flag for it was **false**. We conclude that the lemma holds due to the same arguments employed for the previous cases.

We have considered all possible cases for old and new allocated slots and have shown that they always satisfy the claims. \square

It has just been demonstrated that the collector has full knowledge on which slots have changed since the most recent scan and what were their contents. We now show that the collector can find out what these slots values are in a current cycle as well. These two abilities combined amount to the collector's ability to calculate the asynchronous reference-count of each object, relative to the sliding view of the current cycle.

LEMMA 11.7. *For any object o which is allocated at time $COLLECT_k$ it holds that $o.rc@COLLECT_k = ARC(V_k, o)$.*

PROOF. The claim trivially holds for collection cycle zero, since there are no allocated objects at $COLLECT_0$. To prove that the claim holds for cycle $k > 0$ we assume that it holds for cycle $k - 1$ and that Lemmas 11.8 hold for cycle $k - 1$ and 11.6 hold for cycle k .

It suffices to show that:

- (1) for any slot s due to which rc fields are adjusted by the algorithm the rc field of $V_{k-1}(s)$ is decremented exactly once, during the interval $[COLLECT_{k-1}, COLLECT_k)$, while the rc field of $V_k(s)$ is incremented exactly once during the same interval.
- (2) if $V_{k-1}(s) \neq V_k(s)$ then the algorithm adjusts rc fields due to s .

Consider a memory word s , it is in exactly one of three states, with respect to cycle k : allocated new, allocated, not allocated.

Adjusting rc fields due to allocated new slots. If s has been collected during cycle $k - 1$ then according to Lemma 11.8, the collector decremented the rc field of $V_{k-1}(s)$ when the object containing s was reclaimed. At that point, s assumed the value of **null**, which remained in effect at least until s was reallocated, assuming that Theorem 11.1 holds for cycle $k - 1$.

Another possibility is that the object containing s was reclaimed during cycle m , where $m < k - 1$. Since s is new to cycle k , it was not allocated for cycle $k - 1$ and we have $\sigma_{k-1}(s) \stackrel{\text{def}}{=} HS1_{k-1}$ and by the definition of sliding views we have $V_{k-1}(s) = \mathbf{null}$. Thus, we would expect that no rc field will be decremented due to s . Indeed, since the object containing s was not reclaimed during cycle $k - 1$, no decrement was applied due to s as the result of recursive deletion of cycle $k - 1$. Again, due to Theorem 11.1, we know that when s was reallocated it assumed the value of **null**.

Finally, if s has not been ever allocated before then surely it was not subject to recursive deletion during cycle $k - 1$ and it contained **null** at the time it was allocated.

We conclude that at any rate, by the time s is allocated, it contains **null** and all necessary adjustments have been made to the rc field of $V_{k-1}(s)$ in order to reflect that.

Now we have to show that if $V_k(s) \neq \mathbf{null}$ then the rc field of $V_k(s)$ is incremented and otherwise no rc field is incremented, and, that no rc field is decremented due to s in updating of cycle k .

If no mutator modifies s between its allocation point and before $HS1_k(i)$, then, according to Lemma 11.6, $s \notin Hist_k$ and $\sigma_k(s) \stackrel{\text{def}}{=} HS1_k$. At $\sigma_k(s)$ s still assumes the value of **null** and thus $V_k(s) = \mathbf{null}$. Therefore, we would expect that no rc field will be incremented due to s in cycle k . Since $Hist_k$ does not contain any reference of s , this is actually the case. For the same reason no rc field will be decremented as well.

If, on the other hand, some mutator T_i modifies s between its allocation point and before $HS1_k(i)$ then according to Lemma 11.6, applied for cycle k , $VAL(Hist_k; o) = \{\mathbf{null}\}$. Thus, the collector would adjust rc field due to s during the execution of **Update-Reference-Counters**. No rc field will be decremented due to s as **null** is associated with the slot in $Hist_k$. The collector will then either determine s , or declare it undetermined. If s is determined, it will increment the rc value of the determined value, which we have shown to be equal to $V(s)$. Otherwise, when s is undetermined, the collector adds it to the set $Undetermined_k$. It will subsequently consolidate s during the operation of **Fix-Undetermined-Slots**. The rc field of the resolved value, which also equals $V(s)$, will be incremented exactly once, due to the *Handled* set. No matter whether s is determined or not, we've shown that the rc field of $V_k(s)$ is incremented exactly once.

Adjusting rc fields due to allocated old slots. Since s is not reclaimed during cycle $k - 1$ there is no rc adjustments due to it during the recursive deletion of cycle $k - 1$. It is left to consider the effects due to s in the course of updating during cycle k .

If s is an allocated old slot for cycle k then it may be either modified or non-modified during cycle k .

If s is modified, then (due to Lemma 11.6) $VAL(Hist_k; s) = \{V_{k-1}(s)\}$. Consequently, $V_{k-1}(s).rc$ will be decremented during **Update-Reference-Counters**. Then, s will be either determined or consolidated and the rc value of $V_k(s)$ will be incremented accordingly as shown in the previous paragraphs for new slots.

Otherwise, s is not modified. Then we have $VAL(Hist_k; s) = \emptyset$ and no rc updating due to it occur during cycle k , which is the desired result since $V_{k-1}(s) = V_k(s)$.

Adjusting rc fields due to non-allocated slots. If s has not ever been allocated then the claim trivially holds.

If s has been reclaimed during cycle $k - 1$ then we have shown, while dealing with new slots, that at the time s is reclaimed **null** is assigned to it and the respective rc value of $V_{k-1}(s)$ is decremented accordingly.

Consider a slot s which is not allocated for cycle k that has been most recently

been reclaimed during cycle $m < k - 1$. According to the safety theorem, applied for cycle m , no mutator T_i had access to s after $HS4END_m$. Thus, at $HS1_{k-1}$ no mutator had access to s which leads to $s \notin Hist_k$. Additionally, s could not be the subject of recursive deletion during cycle $k - 1$, because that would have meant that the object containing s was deleted twice in a row, which is contradictory to the safety theorem. We conclude that s is neither the subject of recursive deletion during cycle $k - 1$, nor of rc field updating during cycle k , as desired.

Since we have covered all possible options for the state of s , the claim holds. \square

Building on the foundations provided by the link between the conceptual asynchronous reference-count and the concrete rc field and by the correct implementation of the snooping requirement, proved by Lemma 11.3, we are now ready to prove our main claim.

THEOREM 11.1. *An object o is garbage when it is reclaimed. More specifically, o is not reachable from any mutator T_i after $HS4_k(i)$ and hence o is garbage at $HS4END_k$.*

PROOF. We prove the claim by induction on the cycle number, k . For $k = 0$ we have an empty ZCT_0 therefore no object is reclaimed during this cycle and the claim vacuously holds. For $k > 0$ We prove that the claim is correct provided Lemma 11.7 holds for cycle k .

Let $\{T_1, T_2, \dots, T_n\}$ be the set of all mutators, ordered by the time they respond to the fourth handshake. i.e., $HS4_k(1) < HS4_k(2) < \dots < HS4_k(n)$. Let $\{o_1, \dots, o_m\}$ be the set of objects which **Collect** is invoked for during cycle k , ordered chronologically by the time of the invocation (i.e., o_1 was processed first and o_m —last.)

Consider any object o_j that was processed by **Collect**. We prove that the following invariant holds for o_j :

INVARIANT I1. *For each mutator T_i , o_j was continuously unreachable from T_i in the time interval $[HS4_k(i), HS4_k(n)]$. i.e., was not reachable through any of T_i 's local references and through any global root at any time point in the interval.*

The proof is by double induction: the outer induction variable is j , subscripting the objects that were processed. The inner induction variable is i , denoting the index of mutators in the order they responded to the fourth handshake.

For the basis, we consider o_1 . In order to prove that **I1** holds for o_1 we prove that an additional assertion holds:

INVARIANT I2. $RC(o_1) = 0$ continuously in the time interval $[HS4_k(1), HS4_k(n)]$.

Define **I3** as the logical conjunction of **I1** and **I2**. First we show that **I3** holds for o_1 in the (single-pointed) interval $[HS4_k(1), HS4_k(1)]$. Then we show that given that **I3** holds in the interval $[HS4_k(1), HS4_k(i - 1)]$, then it holds in the interval $[HS4_k(i - 1), HS4_k(i)]$ as well and hence in the entire interval $[HS4_k(1), HS4_k(i)]$.

Invariant **I3**, restricted to the interval $[HS4_k(1), HS4_k(1)]$ simply asserts that o_1 was not directly reachable from any of T_1 's local references and from any global root at $HS4_k(1)$ and that $RC(o_1)@HS4_k(1) = 0$. We prove that this is indeed the case.

Since o_1 was processed the first, **Collect** must have been invoked directly from **Reclaim-Garbage** for it. Thus, $0 = o_1.rc@COLLECT_k$. This implies

$$0 = ARC(V_k, o) \geq RC(o)@HS4_k \implies RC(o)@HS4_k = 0$$

by Lemmas 11.7 and 11.3 and the fact that a reference-count is non-negative. Additionally, o_1 was not directly reachable from T_1 at $HS4_k(1)$, or it would have been marked *local* when T_1 's state was scanned when it responded to the fourth handshake. Finally, o_1 was not directly reachable from any global root at $HS4_k(1)$. To see that this is indeed the case consider any global root r . The collector read r prior to starting the fourth handshake and marked the referenced object *local*. Since the time the collector read r and up to $HS4_k(1)$ all mutators would have marked an object *local* had they stored a reference to the object into r . Thus, at any rate, the object which is referenced by r at $HS1_k$ is marked and thus it cannot be o_1 .

If $n = 1$ then we are done. Otherwise, we prove that **I3** holds for the interval $[HS4_k(i-1), HS4_k(i)]$, where $1 < i \leq n$, provided it holds during the interval $[HS4_k(1), HS4_k(i-1)]$. **I3**, restricted to the interval in question, requires that:

- (1) $RC(o_1) = 0$ continuously during the interval, and
- (2) o_1 was not directly reachable from any of the mutators in the set $P \stackrel{\text{def}}{=} \{T_1, \dots, T_{i-1}\}$ continuously during the interval, and
- (3) o_1 was not directly reachable from any global root continuously during the interval, and
- (4) o_1 was inaccessible from T_i at $HS4_k(i)$.

The inductive hypothesis (on i) assures us that o_1 was not directly reachable from all the mutators in P and from any global root at $HS4_k(i-1)$ and that $RC(o_1)@HS4_k(i-1) = 0$. Examining any possible operation which is scheduled during the interval $[HS4_k(i-1), HS4_k(i)]$ we learn that **I3** remained continuously in effect. We show that any instruction of time $t \in [HS4_k(i-1), HS4_k(i)]$ cannot violate (1),(2) or (3) provided (1),(2) and (3) hold up to time $t-1$ then we show that (4) holds.

- a load cannot violate requirements (1) or (3) simply because it is a load, and not a store. It cannot violate requirement (2) since no object or global root is referring to o_1 , due to the validity of (1) and (3) in previous steps.
- a store operation cannot violate (2) since only a load can.
- a store by a mutator $T_l \in P$ cannot violate (1) or (3) since the operand of the store cannot be o_1 , due to the validity of (2) in previous steps.
- a store by a mutator $T_l \notin P$ cannot violate (1) or (3) because the operand of the store cannot be o_1 since the *Snoop_l* flag is set during the interval and such a step would have marked o_1 *local*.
- to prove that (4) is satisfied: at time $HS4_k(i)$ o_1 is not indirectly reachable, from any mutator or global root, since (1) holds at $HS4_k(i)$. It is not directly reachable from T_i , because that would have caused it being marked *local*. It is not directly reachable from a global root at $HS4_k(i)$ since (3) holds at that moment.

That completes the proof that **I3**, and therefore **I1** in particular, hold for o_1 .

Consider now the object o_j , $1 < j \leq m$. If $o_j.rc@COLLECT_k = 0$ then the same arguments that were employed for o_1 are repeated. Otherwise, we have

$$c \stackrel{\text{def}}{=} o_j.rc@COLLECT_k > 0$$

Since o_j is eventually processed by **Collect** there must have been c slots referencing o_j that were cleared and $o_j.rc$ decremented accordingly, in lines (7-8) of **Collect**. The collector has tested the dirty flags of these slots and found that they were off prior to their processing. Since the dirty flag is off for these slots after $HS4END_k$, no mutator could have changed them after, or at $HS1_k$ and before responding to the fourth handshake (due to Lemma 11.4).

Moreover, since these c slots were contained in objects that were processed prior to o_j the inductive lemma (on objects) apply and we know that no mutator had access to any of the c slots after responding to the fourth handshake. We conclude that these c slots have not been changed after $HS4_k$ and before the collector processed them.

In order to prove **I1** we prove an additional invariant:

INVARIANT I4. *No reference to o_j has been stored during the interval $[HS4_k(1), HS4_k(n)]$ to either a heap slot or a global reference.*

Define **I5** as the logical conjunction of **I1** and **I4**. We prove that **I5** holds for o_j .

We have already said that at $HS4_k(1)$ there existed exactly c references to o_j . All these references were contained in objects that, according to the inductive hypothesis on objects, were unreachable from T_1 at $HS4_k(1)$. Additionally, o_j was not directly reachable from T_1 at $HS4_k(1)$, or it would have been marked local. o_j has not been directly reachable from a global reference at $HS4_k(1)$ since that would have caused it being marked *local*, for the same arguments that were applied for o_1 . Finally, had o_j been indirectly reachable from a global reference r at $HS4_k(1)$ then the chain of references must have passed through some of the c slots which are contained in objects which are assumed to be inaccessible from T_1 at $HS4_k(1)$, contradicting the inductive hypothesis on objects. Thus, **I1**, restricted to the interval $[HS4_k(1), HS4_k(1)]$ holds for o_j .

I4, restricted to the interval $[HS4_k(1), HS4_k(1)]$, holds as well since $HS4_k(1)$ is the time at which T_1 responded to the handshake and naturally it did not execute a store at the same time.

We now show by similar arguments to those applied for o_1 that **I5** restricted to the interval $[HS4_k(i-1), HS4_k(i)]$, where $1 < i \leq n$, holds provided it holds during the interval $[HS4_k(1), HS4_k(i-1)]$. We also use the inductive hypothesis on j that asserts that for any object o_a , $a < j$, **I1** holds for the entire interval $[HS4_k(1), HS4_k(n)]$.

Invariant **I5** applied to o_j and restricted to the interval $[HS4_k(i-1), HS4_k(i)]$ requires that:

- (1) o_j is not reachable continuously during the interval from any local reference of a mutator in P .
- (2) a reference to o_j is not stored during the interval.
- (3) o_j is not reachable continuously during the interval from any global reference.
- (4) o_j is not reachable from T_i at $HS4_k(i)$.

We show that any instruction of time $t \in [HS4_k(i-1), HS4_k(i)]$ cannot violate (1), (2) or (3) provided (1), (2) and (3) hold up to time $t-1$ then we show that (4) holds.

- a load by a mutator T_l could not have made o_j reachable from T_l unless it was reachable from it prior to the load. It also has no effect on the reachability of o_j from other mutators. Therefore such an action cannot violate neither (1) nor (3), assuming (1) and (3) hold for previous steps. Naturally it cannot violate (2).
 - a store by a mutator $T_l \in P$ cannot make o_j reachable for any mutator in P unless o_j has been already reachable from T_l just before the action took place, which is not the case. So a store by T_l preserves (1), (2) and (3) provided (1) and (3) hold for previous steps.
 - a store by a mutator $T_l \notin P$ cannot make o_j reachable from any mutator in P for the following reasons:
 - T_l could not have stored a reference to o_j itself since the $Snoop_l$ flag is set during the interval and such a step would have marked o_j *local*, preventing its processing by **Collect**.
 - T_l could not have stored a reference to x from which o_j is reachable since all references to o_j at the time of the store, by the validity of (2) for previous steps, are a subset of the set of c references that referred to o_j at $HS4_k$. Thus, the chain of references from x to o_j must pass through an object o_a , with $a < j$. The store would have rendered o_a reachable from some mutator in P , which is contradictory to the inductive assumption on o_a .
- So (1), (2) and (3) are not violated by a store by $T_l \notin P$.
- it remains to show that (4) is not violated. Suppose that at $HS4_k(i)$ o_j is reachable from T_i . o_j could not have been directly reachable at the time, or it would have been marked *local*. By the validity of (2) for $HS4_k(i)$ we know that if o_j is reachable from T_i then it is reachable through some object o_a , with $a < j$. This implies that o_a is reachable from T_i at $HS4_k(i)$. Again, a contradiction to the inductive assumption on o_a .

That completes the proof that **I5** and therefore **I1** hold for o_j .

Applying **I1** for any object which is processed we learn that any such object is garbage at $HS4END_k$ (which equals, by definition, $HS4_k(n)$.) Since the objects which are eventually reclaimed are a subset of those processed (the rest have their reclamation deferred to the next cycle) the algorithm is indeed safe. \square

Last but not least we have to prove Lemma 11.8, whose correctness was assumed by Lemma 11.7. The lemma asserts that the collector sensibly de-allocates objects. That is, that it decrements the rc field of slots in a manner which is not discordant with their linkage to the sliding view.

LEMMA 11.8. *Let o be an object which is reclaimed during cycle k and let s be a slot of the object. Then the collector decrements $V_k(s)$ exactly once due to recursive deletion in cycle k .*

PROOF. The claim vacuously holds for cycle $k = 0$. We prove that it holds for cycle $k > 0$ provided Theorem 11.1 and Lemma 11.6 hold for cycle k .

As the reference-count of an object is monotonically non-increasing due to recursive deletion and since an object is processed by **Collect** only when its *rc* field reaches zero, *o* is processed exactly once before being reclaimed.

Since *o* is reclaimed, the collector resets all its slots, including *s*. When the collector considers *s* it probes the value of *Dirty(s)* and finds it off. As noted in Lemma 11.6, *s* could not have been modified by any mutator between responding to the first handshake and fourth handshake. So *s* is not in the digested history for the next cycle.

If $s \notin Hist_k$ then $\sigma_k(s) = HS1_k$. By Lemma 11.5 $Dirty(s) @ \sigma_k(s) = \mathbf{false}$ thus no mutator T_i could have changed *s* between $\sigma_k(s)$ and $HS1_k(i)$. If $s \in Hist_k$ then it must be that $s \in R3_k$. So in that case $\sigma_k(s) = HS2END_k$. At any rate, no mutator T_i changed *s* between $\sigma_k(s)$ and $HS4_k(i)$.

Theorem 11.1 asserts that *s* was inaccessible for any mutator after responding to the fourth handshake.

Assembling these facts we get that any rate *s* was not modified between $\sigma_k(s)$ and the time the collector read its value, prior to resetting it in procedure **Collect**. So the collector indeed decremented the *rc* value of $V_k(s)$. \square

This completes the safety proof of the algorithm.

12. CONCLUSIONS

We have presented a novel on-the-fly reference-counting garbage collector. The new collector is adequate for run on a multiprocessor, in contrast to prior belief that reference-count updates require a large synchronization overhead. The first novelty is in the elimination of a large fraction of the reference-count updates. This elimination results in a drastic reduction in the overhead on reference-count updates, improving the overall throughput of the reference-counting garbage collection. The second novelty is a carefully designed write-barrier that allows concurrent run of the collector with several mutators without requiring synchronized operations in the write barrier. Finally, the proposed collector is on-the-fly, i.e., there is no particular point in which all mutators must be suspended simultaneously. Instead, each mutator cooperates with the collector by being shortly suspended four times during each collection cycle.

We provided a full proof of the collector showing that it is live (eventually reclaims all unreachable objects) and safe (does not reclaim live objects). We have also implemented our collector on Sun's Reference Release 1.2.2 of the Java Virtual Machine and presented measurements demonstrating excellent latency and execution times that are comparable to the original mark-sweep-compact collector. While the pauses were drastically reduced (by two orders of magnitude) the throughput was reduced by at most 8% compared thr original collector.

13. ACKNOWLEDGMENT

We thank Hillel Kolodner for helpful discussions. We thank the anonymous referees and David Wise for their thorough inspection of this paper and for their helpful comments.

REFERENCES

- AHO, A. V., KERNIGHAN, B. W., AND WEINBERGER, P. J. 1988. *The AWK Programming Language*. Addison-Wesley, Reading, MA.
- APPEL, A. W., ELLIS, J. R., AND LI, K. 1988. Real-time concurrent collection on stock multiprocessors. *ACM SIGPLAN Not.* 23, 7, 11–20.
- AZATCHI, H., LEVANOI, Y., PAZ, H., AND PETRANK, E. 2003. An on-the-fly mark and sweep garbage collector based on sliding views. See OOPSLA [2003], 269–281.
- AZATCHI, H. AND PETRANK, E. 2003. Integrating generations with advanced reference counting garbage collectors. In *International Conference on Compiler Construction (CC'2003)*. Lecture Notes in Computer Science, vol. 2622. Springer, Berlin, 185–199.
- BACON, D., ATTANASIO, D., LEE, H., AND SMITH, S. 2001. Java without the coffee breaks: A nonintrusive multiprocessor garbage collector. In *Proceedings of SIGPLAN 2001 Conference on Programming Languages Design and Implementation*. ACM SIGPLAN Not. ACM Press, New York, 92–103.
- BACON, D. AND RAJAN, V. 2001. Concurrent cycle collection in reference counted systems. In *Proceedings of 15th European Conference on Object-Oriented Programming, ECOOP 2001*. Budapest.
- BAKER, H. G. 1978. List processing in real-time on a serial computer. *Commun. ACM* 21, 4, 280–94.
- BAKER, H. G. 1994. Minimising reference count updating with deferred and anchored pointers for functional data structures. *ACM SIGPLAN Not.* 29, 9.
- BARTH, J. M. 1977. Shifting garbage collection overhead to compile time. *Commun. ACM* 20, 7 (July), 513–518.
- BLACKBURN, S. AND MCKINLEY, K. 2003. Ulterior reference counting: Fast garbage collection without a long wait. See OOPSLA [2003], 344–358.
- BOEHM, H.-J., DEMERS, A. J., AND SHENKER, S. 1991. Mostly parallel garbage collection. *ACM SIGPLAN Not.* 26, 6, 157–164.
- CHIKAYAMA, T. AND KIMURA, Y. 1987. Multiple reference management in Flat GHC. In *4th International Conference on Logic Programming*. MIT Press, 276–293.
- COLLINS, G. E. 1960. A method for overlapping and erasure of lists. *Commun. ACM* 3, 12 (Dec.), 655–657.
- CRAMMOND, J. 1988. A garbage collection algorithm for shared memory parallel processors. *International Journal Of Parallel Programming* 17, 6, 497–522.
- DETREVILLE, J. 1990. Experience with concurrent garbage collectors for Modula-2+. Tech. Rep. 64, DEC Systems Research Center, Palo Alto, CA. Aug.
- DEUTSCH, L. P. AND BOBROW, D. G. 1976. An efficient incremental automatic garbage collector. *Commun. ACM* 19, 9 (Sept.), 522–526.
- DIJKSTRA, E. W., LAMPORT, L., MARTIN, A. J., SCHOLTEN, C. S., AND STEFFENS, E. F. M. 1978. On-the-fly garbage collection: An exercise in cooperation. *Commun. ACM* 21, 11 (Nov.), 965–975.
- DOLIGEZ, D. AND GONTHIER, G. 1994. Portable, unobtrusive garbage collection for multiprocessor systems. In *Conference Record of the Twenty-first Annual ACM Symposium on Principles of Programming Languages*. ACM SIGPLAN Not. ACM Press, New York, 70–83.
- DOLIGEZ, D. AND LEROY, X. 1993. A concurrent generational garbage collector for a multi-threaded implementation of ML. In *Conference Record of the Twentieth Annual ACM Symposium on Principles of Programming Languages*. ACM SIGPLAN Not. ACM Press, New York, 113–123.
- DOMANI, T., KOLODNER, E. K., LEWIS, E., SALANT, E. E., BARABASH, K., LAHAN, I., PETRANK, E., YANOVER, I., AND LEVANOI, Y. 2000. Implementing an on-the-fly garbage collector for Java. See Hosking [2000], 155 – 166.
- ENDO, T., TAURA, K., AND YONEZAWA, A. 1997. A scalable mark-sweep garbage collector on large-scale shared-memory machines. In *Proceedings of the Conference on High Performance Networking and Computing (SC'97)*. ACM Press, New York, 1 – 14.

- FLOOD, C., DETLEFS, D., SHAVIT, N., AND ZHANG, C. 2001. Parallel garbage collection for shared memory multiprocessors. In *Usenix Java Virtual Machine Research and Technology Symposium (JVM '01)*. Monterey, CA.
- FURUSOU, S., MATSUOKA, S., AND YONEZAWA, A. 1991. Parallel conservative garbage collection with fast allocation. In *OOPSLA/ECOOP '91 Workshop on Garbage Collection in Object-Oriented Systems, Addendum to OOPSLA'91 Proceedings*, P. R. Wilson and B. Hayes, Eds.
- GOLDBERG, A. AND ROBSON, D. 1983. *Smalltalk-80: The Language and its Implementation*. Addison-Wesley, Reading, MA.
- HALSTEAD, R. H. 1985. Multilisp: A language for concurrent symbolic computation. *ACM Trans. Program. Lang. Syst.* 7, 4 (Oct.), 501–538.
- HERLIHY, M. AND MOSS, J. E. B. 1990. Non-blocking garbage collection for multiprocessors. Tech. Rep. CRL 90/9, DEC Cambridge Research Laboratory, Cambridge, MA.
- HOSKING, A. L., MOSS, J. E. B., AND STEFANOVIĆ, D. 1992. A comparative performance evaluation of write barrier implementations. In *OOPSLA'92 ACM Conference on Object-Oriented Systems, Languages and Applications*, A. Paepcke, Ed. ACM SIGPLAN Not., vol. 27(10). ACM Press, New York, 92–109.
- HOSKING, T., Ed. 2000. *ISMM 2000 Proceedings of the Second International Symposium on Memory Management*. ACM SIGPLAN Not., vol. 36(1). ACM Press, New York.
- HUDSON, R. L. AND MOSS, J. E. B. 2003. Copying garbage collection without stopping the world. *Concurrency and Computation: Practice and Experience* 15, 3-5, 223–261.
- JONES, R. E. AND LINS, R. 1996. *Garbage Collection: Algorithms for Automatic Dynamic Memory Management*. Wiley.
- KOLODNER, E. K. AND PETRANK, E. 2004. Parallel copying garbage collection using delayed allocation. *Parallel Processing Letters* 14, 2. To appear.
- LEVANONI, Y. AND PETRANK, E. 2001. n on-the-fly reference counting garbage collector for Java. In *OOPSLA'01 ACM Conference on Object-Oriented Systems, Languages and Applications*. ACM SIGPLAN Not., vol. 36(10). ACM Press, New York.
- LINS, R. D. AND VASQUES, M. A. 1991. A comparative study of algorithms for cyclic reference counting. Tech. Rep. 92, Computing Laboratory, The University of Kent at Canterbury. Aug.
- MCBETH, J. H. 1963. On the reference counter method. *Commun. ACM* 6, 9, 575.
- MILLER, J. S. AND EPSTEIN, B. 1990. Garbage collection in MultiScheme. In *US/Japan Workshop on Parallel Lisp, LNCS 441*. Springer, Berlin, 138–160.
- OOPSLA 2003. *OOPSLA'03 ACM Conference on Object-Oriented Systems, Languages and Applications*. ACM Press, New York.
- OSSIA, Y., BEN-YITZHAK, O., GOFT, I., KOLODNER, E. K., LEIKEHMAN, V., AND OWSHANKO, A. 2002. A parallel, incremental and concurrent GC for servers. In *Proceedings of SIGPLAN 2002 Conference on Programming Languages Design and Implementation*. ACM Press, New York, 129–140.
- O'TOOLE, J. W. AND NETTLES, S. M. 1994. Concurrent replicating garbage collection. In *Proceedings of the 1994 ACM conference on LISP and functional programming*. ACM Press, New York, 34 – 42.
- PARK, Y. G. AND GOLDBERG, B. 1995. Static analysis for optimising reference counting. *Inf. Process. Lett.* 55, 4 (Aug.), 229–234.
- PLAKAL, M. AND FISCHER, C. N. 2000. Concurrent garbage collection using program slices on multithreaded processors. See Hosking [2000].
- PRINTEZIS, T. AND DETLEFS, D. 2000. A generational mostly-concurrent garbage collector. See Hosking [2000], 143 – 154.
- RIANY, Y., SHAVIT, N., AND TOUITOU, D. 1995. Towards a practical snapshot algorithm. In *Proceedings of the 3rd Israeli Symposium on the Theory of Computing and Systems*. IEEE Press, 58–173.
- ROTH, D. J. AND WISE, D. S. 1998. One-bit counts between unique and sticky. In *ISMM'98 Proceedings of the First International Symposium on Memory Management*, R. Jones, Ed. ACM SIGPLAN Not., vol. 34(3). ACM Press, New York, 49–56.

- SOBALVARRO, P. 1988. A lifetime-based garbage collector for Lisp systems on general-purpose computers. Tech. Rep. AITR-1417, MIT AI Lab, Cambridge, MA. Feb. Bachelor of Science thesis.
- SPEC BENCHMARKS. 1998,2000. Standard Performance Evaluation Corporation. <http://www.spec.org/>.
- STEELE, G. L. 1975. Multiprocessing compactifying garbage collection. *Commun. ACM* 18, 9 (Sept.), 495–508.
- STOYE, W. R., CLARKE, T. J. W., AND NORMAN, A. C. 1984. Some practical methods for rapid combinator reduction. In *Conference Record of the 1984 ACM Symposium on Lisp and Functional Programming*, G. L. Steele, Ed. ACM Press, New York, 159–166.
- WALL, L. AND SCHWARTZ, R. L. 1991. *Programming Perl*. O'Reilly and Associates, Inc.
- WEIZENBAUM, J. 1963. Symmetric list processor. *Commun. ACM* 6, 9 (Sept.), 524–544.
- WISE, D. S. 1993. Stop and one-bit reference counting. *Inf. Process. Lett.* 46, 5 (July), 243–249.
- YUASA, T. 1990. Real-time garbage collection on general-purpose machines. *Journal of Software and Systems* 11, 3, 181–198.
- ZORN, B. 1990. Barrier methods for garbage collection. Tech. Rep. CU-CS-494-90, University of Colorado, Boulder. Nov.