

# **NBR: Neutralization Based Reclamation**

Ajay Singh University of Waterloo Canada ajay.singh1@uwaterloo.ca Trevor Brown University of Waterloo Canada trevor.brown@uwaterloo.ca Ali Mashtizadeh University of Waterloo Canada mashti@uwaterloo.ca

## Abstract

*Safe memory reclamation* (SMR) algorithms suffer from a trade-off between bounding unreclaimed memory and the speed of reclamation. Hazard pointer (HP) based algorithms bound unreclaimed memory at all times, but tend to be slower than other approaches. Epoch based reclamation (EBR) algorithms are faster, but do not bound memory reclamation. Other algorithms follow hybrid approaches, requiring special compiler or hardware support, changes to record layouts, and/or extensive code changes. Not all SMR algorithms can be used to reclaim memory for all data structures.

We propose a new neutralization based reclamation (NBR) algorithm that is often faster than the best known EBR algorithms and achieves bounded unreclaimed memory. It is non-blocking when used with a non-blocking operating system (OS) kernel, and only requires atomic read, write and CAS. NBR is straightforward to use with many different data structures, and in most cases, requires similar reasoning and programmer effort to two-phased locking. NBR is implemented using OS signals and a lightweight handshaking mechanism between participating threads to determine when it is safe to reclaim a record. Experiments on a lockbased binary search tree and a lazy linked list show that NBR significantly outperforms many state of the art reclamation algorithms. In the tree, NBR is faster than next best algorithm, DEBRA, by up to 38% and HP by up to 17%. And, in the list, NBR is 15% and 243% faster than DEBRA and HP, respectively.

# $\label{eq:ccs} \textit{CCS Concepts:} \bullet \textit{Computing methodologies} \rightarrow \textit{Concurrent computing methodologies}.$

Keywords: safe memory reclamation, concurrency.

#### **ACM Reference Format:**

Ajay Singh, Trevor Brown, and Ali Mashtizadeh. 2021. NBR: Neutralization Based Reclamation. In 26th ACM SIGPLAN Symposium

PPoPP '21, February 27-March 3, 2021, Virtual Event, Republic of Korea © 2021 Copyright held by the owner/author(s). ACM ISBN 978-1-4503-8294-6/21/02. https://doi.org/10.1145/3437801.3441625 on Principles and Practice of Parallel Programming (PPoPP '21), February 27-March 3, 2021, Virtual Event, Republic of Korea. ACM, New York, NY, USA, 16 pages. https://doi.org/10.1145/3437801.3441625

## 1 Introduction

Fundamentally, *safe memory reclamation* (SMR) is about answering the question: When is it safe to free a record? Unlike garbage collection, which is automatic, SMR requires a program to invoke a *retire* operation on each record at some point after it becomes *garbage* (i.e., is *unlinked* from the data structure). The task of an SMR algorithm is to eventually *free* an unlinked record once no thread holds a pointer to it [7, 13, 35].

The challenge of SMR in concurrent data structures comes from use-after-free conflicts between threads, where one thread accesses a record that is concurrently freed by another. For example, consider a lazy-list where one thread is searching and another is deleting. The first thread obtains a reference to a record and stores it in a local variable. The other thread unlinks and frees. At this point the first thread's reference is no longer safe as the record it points to has been freed.

Researchers have developed a rich variety of SMR algorithms, with a diverse spectrum of desirable properties, idiosyncrasies and limitations. After experimenting with SMR algorithms and observing the state of art [2-7, 13-16, 18, 19, 21, 27, 29, 32, 35, 37, 42, 46], we identified the following set of desirable properties. [P1] Performance: reclamation operations should ideally offer both low latency and high throughput. [P2] Bounded Garbage: The number of records that are unlinked but not yet reclaimed should be bounded, even if threads experience halting failures or long delays. [P3] Usability: Intrusive changes to code, data, and the build environment, should be minimized. [P4] Consistency: Performance should not be drastically affected by changes in the workload (e.g., when shifting between read-intensive and update-intensive workloads). Additionally, there should be minimal performance degradation when the system is oversubscribed (with more threads than cores). [P5] Applicability: The algorithm should be usable with as many data structures as possible.

To set the stage for our contribution, we must first discuss other approaches. We broadly categorize existing work into: hazard pointer-based reclamation (HPBR), quiescent statebased reclamation (QSBR), epoch-based reclamation (EBR), reference counting based reclamation (RCBR), and hybrid algorithms that combine the aforementioned approaches [29].

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for thirdparty components of this work must be honored. For all other uses, contact the owner/author(s).

In general, QSBR and EBR are fast but do not bound garbage, HPBR has bounded garbage but is not fast, and RCBR is neither fast nor does it bound garbage (in case retired nodes have pointer cycles [18]). Hybrid approaches have generally focused on achieving P1 and P2 simultaneously, usually by combining EBR (for its speed) with some variant of HPBR (to bound garbage), with varying levels of success.

The hybrid algorithm that most closely resembles our approach is DEBRA+ [13], a variant of EBR (with a restricted form of HPBR) that is designed for lock-free data structures. DEBRA+ is fast, and it achieves bounded garbage via a neutralizing mechanism based on POSIX signals and data structure specific recovery code. A thread whose reclamation is delayed by a slow thread will send a *neutralizing signal* to the slow thread. Upon receipt of a neutralizing signal, a thread executes its recovery code and then restarts its data structure operation, allowing reclamation to continue, ultimately guaranteeing a bound on the number of unreclaimed records. However, this bound on garbage comes at the cost of both usability and applicability, as users need to write data structure specific recovery code that is not always straightforward, or even possible. Moreover, it is not clear how DEBRA+ could be used for lock-based data structures, since neutralizing a thread that holds a lock could cause deadlock.

**Contribution**: Existing SMR algorithms all have significant shortcomings in their attempts at satisfying properties P1 through P5. This motivated us to propose a new *Neutralization Based Reclamation* algorithm (*NBR*) that matches or outperforms existing SMR algorithms [P1], bounds garbage [P2], is simple to use [P3], exhibits consistent performance, even on oversubscribed systems [P4], and is applicable to a large class of data structures, some of which are not supported by popular SMR algorithms [P5].

NBR's neutralization technique is similar to that of DE-BRA+, with a few key differences. In NBR, each thread places unlinked objects in a thread-local buffer, and when the buffer's size exceeds a predetermined threshold, the thread sends a neutralizing signal to all other threads. Upon receipt of such a signal, a thread checks whether its current data structure operation has already done any writes to shared memory, and if not, restarts its operation (using the C/C++ procedures sigset jmp and siglong jmp). Otherwise, it finishes executing its operation. In contrast, to guarantee bounded garbage in DEBRA+, a thread must restart even if it has already written to shared memory-a design decision that limits DEBRA+'s applicability to specific lock-free data structures, and necessitates data structure specific recovery code. *NBR* does not require any recovery code, and can be used with nearly all structures that DEBRA+ supports and many others structures DEBRA+ does not, including some lockbased algorithms, such as a lock-based binary search tree with lock-free searches [17] (DGT).

We also present an optimized version of *NBR* called *NBR+* in which threads send fewer signals, and yet reclaim memory

Ajay Singh, Trevor Brown, and Ali Mashtizadeh

more often. This is accomplished by allowing threads to infer when memory can be freed simply by passively observing the signals sent in the system. Finally, as our experiments show, NBR+ is highly efficient, significantly outperforming the state of the art in SMR in various data structure workloads on a large-scale 4-socket Intel system.

The rest of the paper is structured as follows. Related work is surveyed in Section 2. In Section 3, we introduce the model. Section 4 describes our basic algorithm *NBR*, and characterizes its applicability. We describe an optimized version *NBR+* in Section 5, followed by a brief discussion on *NBR*'s correctness in Section 5.4. Finally, experiments appear in Section 6, followed by conclusions in Section 7.

## 2 Related Work

Although detailed surveys of *safe memory reclamation* already exist in earlier works [13, 29], we would like to study existing techniques specifically through the lens of the desirable properties defined above.

**RCBR** involves explicitly counting the number of incoming pointers to a record, and typically storing this count alongside the record. The inclusion of this metadata in records complicates any advanced pointer arithmetic or implicit pointers, and can require changes to record layouts (or the use of a custom allocator) as well as size. RCBR typically requires a programmer to invoke a *deref* operation to dereference a pointer (and sometimes to explicitly invoke operations for read, write and CAS) [6, 18, 32, 37], adding significant overhead and programmer effort [opposing P1, P3]. Programmer intervention is also needed to identify and break pointer *cycles* in garbage records.

HPBR incurs significant overhead every time a new record is encountered, as a thread must first announce a hazard pointer (HP) to it in a shared location, then issue a memory fence (or use an atomic exchange instruction to announce the hazard pointer) and then check whether the record has already been unlinked [19, 32, 35] [opposing P1, P3]. If the record has been unlinked, the data structure operation trying to access it must be *restarted* (a data structure specific action). Correctly dealing with such failure cases can require extensive code changes. This may also require the programmer to reprove the data structure's progress guarantees [13]. Additionally, it is not clear how HPs could be used with data structures that allow threads to traverse pointers in unlinked records [13], and there are many examples of such data structures, e.g., [1, 8, 10, 20, 23, 25, 31, 36, 39, 43] [opposing P5]. (In such data structures, a search can potentially pass through many unlinked records, and yet end up back in the data structure, at the appropriate location.)

The latter limitation was addressed by *Beware and Cleanup* a hybrid of RCBR and HPBR [27]. However, this algorithm requires a programmer to write a data structure specific *cleanup* procedure that changes all pointers in an unlinked record to point to *current* records *in the data structure*. This cleanup code ensures bounded garbage for data structures that allow traversing unlinked records, but the algorithm has higher overhead than either RCBR or HPBR and requires significant programmer effort [opposing P1, P3].

**QSBR and EBR** [26, 29, 34] both use the observation that, in many data structures, threads do *not* carry pointers obtained in one data structure operation forward for use in a subsequent operation. QSBR and EBR each have a simple interface in which the programmer need only invoke a specific operation at the start and end of a data structure operation. Unlike the approaches above, QSBR and EBR avoid all perrecord and per-access overheads. A thread can reclaim its garbage records whenever it detects that all other threads have started a new data structure operation (and hence *forgotten* all pointers to said garbage records). However, in the event that a thread halts or is delayed, the amount of unreclaimed garbage can grow unboundedly [opposing P2].

DEBRA+ (described above) was introduced by Brown in 2015 [13]. In the same paper, an algorithm called DEBRA was proposed which, to the best of our knowledge, is the fastest EBR algorithm. DEBRA does not bound the number of unreclaimed records (garbage), but was shown to be faster than DEBRA+. Note that our experiments show NBR+ often *outperforms* DEBRA.

Since DEBRA, numerous hybrid algorithms offering bounded garbage have appeared, for example, Hyaline (HY) [37], Hazard Eras (HE) [42], Interval Based Reclamation (IBR) [46], and Wait Free Eras (WFE)[38]. All of these algorithms use per-record metadata to encode the times at which a record is allocated and unlinked, and the code instrumentation needed is similar to HPs [opposing P3]. HPs require per-record reserve and unreserve calls and fallback code (to restart the operation) in case the reservation fails (because the record is, or might be, unlinked). It is unclear how HE, IBR and WFE could be used with data structures that allow traversing unlinked records [opposing P5]. As we will see in our experiments, these algorithms also incur non-trivial overhead [opposing P1].

Various other algorithms utilize operating system features such as forced context switches [5], POSIX signals [3, 4], and hardware transactional memory [2, 21]. OSense [5] is a hybrid algorithm that uses QSBR as a *fast* code path, and HPs with forced context switches as a *slow* code path to bound garbage. However, in the event of long thread delays, reclamation can only proceed on the slow path, which is as slow as HPBR. QSense has been shown to be slower than EBR [5] [opposing P1]. None of [2, 4, 5] can be used with data structures that allow threads to traverse unlinked records [opposing P5]. Forkscan (FS) [3] was succeeded by Thread-Scan (TS) [4], which addressed this issue, but FS assumes the programmer will not use advanced pointer arithmetic techniques (or implicit pointers) [opposing P3]. Additionally, FS has been shown to be slower than HPs in several workloads [3] [opposing P1].



**Figure 1.** Visualizing the form of a *data-structure* operation for which *NBR* can be used. The thread performing this operation can be neutralized in the read-phase ( $\Phi_{read}$ ). However, it cannot be neutralized in the write-phase ( $\Phi_{write}$ ). end $\Phi_{read}$  marks the beginning of the operation's  $\Phi_{write}$ .

Optimistic Access (OA) and Automatic Optimistic Access (AOA) [15, 16] proposed a particularly interesting approach: they optimistically allow threads to accesses reclaimed nodes, and verify after the fact that the access was safe. This requires an assumption that either (a) memory will not be freed to the OS, or (B) any resulting trap/exception (such as a segmentation fault) will be caught and handled [opposing P3]. Additionally, OA requires the programmer to transform data structures into a normalized form [45] (which is similar to, but not the same as, the form we assume in this paper), and instrument every read/write/CAS [opposing P3]. AOA automates this transformation with compiler support (for data structures that *have* a normalized form). Unfortunately, it doesn't appear that AOA has been ported to modern compilers. The need for a normalized form was eliminated in Free Access (FA) [14], which used a compiler extension to perform automatic instrumentation of writes and blocks of consecutive independent reads. FA is a general technique that has been shown to have comparable performance to HPBR [14] [opposing P1]. In contrast, our work targets applications that can benefit from the high performance handcrafted SMR.

## 3 Model

We consider an *n* thread asynchronous shared memory system. Threads can perform atomic read, write, compare-andswap (CAS) and fetch-and-add (FAA). A data structure consists of a set of records which are accessible from a *root* (e.g., the head of a list). A *record* can be viewed as a set of fields. Each record can be in one of five states throughout its lifecycle: (1) *allocated*: record allocated from heap but not accessible through the root, (2) *reachable*: the record can be reached by following references from the root, (3) *unlinked*: is not reachable from (any) root but threads may still have references to it in thread private memory, (4) *safe*: a record is unlinked and no thread has a reference to it, and (5) *reclaimed* (or freed): a record is returned to the OS. In states 3 and 4, a record is *garbage*.

## 4 NBR

### 4.1 Assumptions on the data structure

*NBR* requires a data structure's operations to have (or be restructured into) the following form. This form is needed for the neutralization mechanism wherein an operation that

has not yet written to shared memory is forced to restart. The required form is described as a sequence of *phases* (illustrated in Figure 1), with specific rules in each phase. Along the way, we discuss potential pitfalls for readers unfamiliar with the use of sigsetjmp and siglongjmp.

**Phase 0:** preamble. Accesses (reads/writes/CASs) to global variables are permitted. System calls (heap allocation/deallocation, file I/O, network I/O, etc.) are permitted. Access to shared records, for example, nodes of a shared data structure, is *not* permitted.

**Phase 1:**  $\Phi_{read}$  (read phase). Reading global variables is permitted and reading shared records is permitted if pointers to them were obtained during this phase (e.g., by traversing a sequence of shared objects by following pointers starting from a global variable—i.e., a *root*). Writes/CASs to shared records, writes/CASs to shared globals, and system calls, are *not* permitted.

To understand the latter restriction, suppose an operation allocates a node using malloc during its  $\Phi_{read}$ , and before it uses the node, the thread performing the operation is neutralized. This would cause a memory leak.

Additionally, writes to thread local data structures are not recommended. To see why, suppose a thread maintains a thread local doubly-linked list, and also updates this list as part of the  $\Phi_{read}$  of some operation on the shared data structure. If the thread is neutralized in middle of its update to this local list, it might corrupt the structure of the list. <sup>a</sup> **Phase 2:** reservation. This is a *conceptual stage* that does not necessarily correspond to any data structure code. However, this is where a key *NBR* operation will be invoked. At this point, one must be able to identify all shared objects that will be modified by the operation in the next phase, so they can be provided to *NBR*. We call these *reserved* records.

**Phase 3:**  $\Phi_{write}$  (write phase). Accesses (reads/writes/CASs) to global variables, and system calls, are permitted. Accesses (including write/CAS) to shared records are permitted only if the records are *reserved*. To understand what could go wrong if this restriction is violated, we need to better understand *NBR*, so we will return to this restriction with an example in Section 4.4.

Finally, threads not executing a data structure operation are said to be in a *quiescent phase* (essentially the same as phase 0).

#### 4.2 Overview of NBR

In *NBR*, each thread accumulates records that it has unlinked in a private buffer (or *limbo bag*). When the size of a thread *T*'s buffer exceeds a predetermined threshold, the thread sends a neutralizing signal to all other threads. Upon receipt of such a signal, the behaviour of a thread T' depends on which phase it is executing in.

If T' is in a quiescent phase, or preamble (Phase 0), it holds no pointers to shared records, and does not prevent T from reclaiming records in its buffer. T' simply continues executing (effectively ignoring the signal).

On the other hand, if T' is in  $\Phi_{read}$ , it may hold pointers to records in T's buffer. If T' were to continue executing, it would have to prevent T from reclaiming records. Note, however, that T' has not yet performed any modifications to any shared records (since it is still in  $\Phi_{read}$ ). So, T' can simply discard all of its pointers (that are in its private memory), and jump back to the start of its  $\Phi_{read}$ , without leaving any shared data structures in an inconsistent state. To implement this jump, every data structure operation invokes sigsetjmp at the start of its  $\Phi_{read}$ , which creates a *checkpoint* (saving the values of all stack variables). A thread can subsequently invoke siglong jmp to return to the last place it performed sigsetjmp (and restore the values of all stack variables).<sup>b</sup> It can then retry executing its  $\Phi_{read}$ , traversing a new sequence of records, starting from the root, without any risk of accessing any records freed by T (since those are no longer reachable).

The subtlety in *NBR* arises when T' is in  $\Phi_{write}$ . In this case, T' may hold pointers to records in the buffer of T. Thus, if it continues executing, T' must prevent T from reclaiming these records. Moreover, since T' may have modified some shared records (but not completed its operation yet), we cannot simply restart its data structure operation, or we may leave the data structure in an inconsistent state. So, T' will *not* restart its operation. Instead, it will simply *continue executing* wherever it was when it received the signal (effectively ignoring the signal). At this point, the reader might wonder how we *simultaneously avoid*:

- **a.** Blocking the reclamation of *T*, and
- **b.** The possibility that *T*′ continues executing and it accesses a record freed by *T*.

The solution lies in the *reservation* phase (Phase 2) of T'. During the reservation phase of T', just before it begins its  $\Phi_{write}$ , T' *reserves* all of the shared records it will access in its  $\Phi_{write}$  by *announcing* pointers to them in a shared array. These reservations serve a similar purpose to hazard pointers, but are quite different from HP in terms of performance and safety guarantees. This is discussed further in Section 5.3. By the time T' is in its  $\Phi_{write}$  (so it will ignore any neutralization signals), its reservations are visible to T, and T can refer to these reservations to avoid reclaiming any of those records.

In short, operations in the  $\Phi_{read}$  discard their pointers and restart, and operations in the  $\Phi_{write}$  must have reserved them. This empowers the *reclaimers* to assume that *readers* lose all of their pointers in response to neutralizing, and the

<sup>&</sup>lt;sup>a</sup>In some cases it is safe to write to thread local storage (TLS). For example, a thread could use TLS to maintain statistics that are *supposed to* persist despite neutralization. Similarly, some idempotent or atomic changes to TLS should remain correct even if one is neutralized. Proceed with caution.

<sup>&</sup>lt;sup>b</sup>Technically sigsetjmp saves only the current stack *frame*. Stack variables defined deeper on the stack will not necessarily be saved or restored.

Algorithm 1 NBR. Assumes, max number of reservations are less than the limboBag size.

#### thread local variable:

```
1: int tid
                                                     ▷ current thread id
```

- 2: record \*limboBag ▷ per-thread list of unlinked records. Maxsize:S
- 3: bool restartable ▷ local var to track  $\Phi_{read} / \Phi_{write}$
- 4: record \*tail ▶ Pointer to last record in limboBag

#### shared variable:

5: atomic<record\*> reservations[N][R] ▷ N:#threads, R:max reserved records.  $|\mathbf{R}| << |\mathbf{S}|$ .

```
6: procedure BEGIN\Phi_{read}()
```

```
reservations[tid].clear();
7:
      CAS(&restartable, 0, 1);
8:
   end procedure
9:
   procedure END\Phi_{read}(\{rec_1, rec_2 \cdots rec_R\})
10:
      reservations[tid] = {rec_1, rec_2 \cdots rec_R};
11:
      CAS(&restartable, 1, 0);
12:
13: end procedure
   procedure RETIRE(rec)
14:
      if isLimboBagTooLarge() then
15:
          signalAll( );
16:
          reclaimFreeable(tail);
17:
18:
      end if
      limboBag[tid].append(rec);
19:
20: end procedure
   procedure RECLAIMFREEABLE(tail)
21:
      A = collectReservations( );
22:
      R = limboBag[tid].remove(A, tail);
23:
      free({R});
24:
   end procedure
25:
```

writers lose all pointers that are not reserved. As a result, once a thread sends a neutralizing signal to all other threads, it can scan all reservations, and free any records in its buffer (limbo bag) that are not reserved.

#### 4.3 Implementation of NBR

Algorithm 1 shows the pseudocode for NBR. Each thread collects unlinked records in its limboBag (line 2), and maintains a local restartable variable that indicates whether the thread should jump back to the start of its  $\Phi_{read}$  in the event that it receives a neutralization signal (line 3). We say the thread is restartable if restartable is true (1), and non-restartable otherwise. Additionally, each thread, before entering the  $\Phi_{write}$ , reserves all records it might access in a single-writer multireader (SWMR) reservation array, (line 5). We assume the maximum number *R* of reserved records is strictly less than maximum size *S* of a limbo bag.

A thread in  $\Phi_{read}$  clears its reservations (if any), and then changes restartable to true using a CAS (Line 8). This CAS might initially seem strange, since it is performed on a singlewriter variable and cannot fail. The CAS prevents instruction

reordering on x86-64 architectures (additional fences may be needed for more relaxed memory models). More specifically, the goal of CAS at line 8 is to ensure that a thread T becomes restartable before any subsequent reads of shared records. If this CAS were simply an atomic write (rather than a readmodify-write instruction), it would be possible for T's reads of shared records to be reordered before this write. In other words some reads of shared records in  $\Phi_{read}$  may appear to occur in *preamble* (or previous  $\Phi_{write}$ ) due to instruction reordering. This could end up breaking the rule that says access to shared records is not permitted in preamble (phase 0) as discussed in Section 4.1. As a result, the thread, which is not yet restartable, might ignore a neutralization signal and access a freed record.

Just before a thread *T* enters a  $\Phi_{write}$ , it announces a set of reservations, and then changes restartable to false using CAS (Line 12). This CAS is used to broadcast the reservations to other threads. More specifically, a CAS by thread T at line 12 implies a memory fence, which ensures that all of the reservations (announced at the previous line 11) are visible to other threads before T changes restartable to false. If this CAS were a simple write, it would be possible for a reclaimer to miss some reservations of *T*, and erroneously free those records<sup>c</sup>

In other words, the following incorrect execution may occur on x86/64 if a write is used instead of CAS: a thread T reserves record *rec* and writes 0 to restartable. Suppose the reservations of thread T remain in the processor's store buffer, and are not visible to other threads yet. Then, another thread T' sends a neutralizing signal to T, scans the reservations and does not see rec, and consequently frees rec. Upon receiving the signal, T' will not restart since it has already written 0 to restartable.<sup>d</sup> Instead, it continues executing, and dereferences rec (accessing a freed record).

The retire operation (line 14) begins by checking whether the size of the limbo bag is above a predetermined threshold (32k in our experiments), at line 15. If so, it sends a neutralizing signal to all threads using signalAll (line 16), and then proceeds to reclaim all safe (i.e., unreserved) records (line 17). Otherwise, it simply adds rec to limboBag.

The reclaimFreeable procedure frees all records (up to the last record pointed to by thread local pointer, tail) in the limboBag that are not reserved (line 21). It first scans reservations array of all other threads and collects the reserved records in set A (line 22). Then it removes the retired records, which are not in A (set of reserved records), up to the *tail* of the *limboBag* using remove(A, tail) at line 23. Finally, it frees the *safe* set of records *R* at line 24.

<sup>&</sup>lt;sup>c</sup>Instead of using CAS, on modern x86/64 machines we can use the more efficient xchg instruction. See Section 11.5.1 of https://www.amd.com/system/ files/TechDocs/47414\_15h\_sw\_opt\_guide.pdf for further details.

<sup>&</sup>lt;sup>d</sup>If *T* and *T'* are executing on different processors, then *T* will not see the effects of any pending writes in the store buffer of T', but T' will see the effects its own pending writes in order to maintain sequential consistency.

After discussing the implementation of *NBR* we can now elaborate on how *readers*, *writers* and *reclaimers* collaborate to achieve safe memory reclamation.

**4.3.1 Reader-reclaimer handshake** Each thread T' at the time of BEGIN $\Phi_{read}$  saves its execution state (program counter and stack frame) using *sigsetjmp* so that when it becomes restartable it can jump back to this state upon receiving a neutralizing signal. When a *reclaimer* T sends a neutralization signal to thread T', the operating system causes the control flow of T' to be interrupted, so that T' will immediately execute a *signal handler* if T' is currently running. (Otherwise, if T' is not currently running, the next time it is scheduled to run it will execute the signal handler before any other steps.) The signal handler determines whether T' is restartable by reading the local *restartable* variable. If the thread is restartable, then the signal handler will invoke siglongjmp and jump back to the start of the  $\Phi_{read}$  (so it is as if T' never started the  $\Phi_{read}$ ).

This behaviour represents a sort of two-step *handshake* between *readers* (threads in  $\Phi_{read}$ ) and *reclaimers* (threads executing lines 16 and 17 in retire) to avoid scenarios where a reader might access a freed record. A *reclaimer* guarantees that before *reclaiming* any of its *unlinked* records it will signal all threads, and all *readers* guarantee that they will relinquish any reference to unsafe records when they receive a neutralization signal.

**4.3.2** Writers handshake (1) Each *reclaimer* signals all threads before starting to reclaim any records. When a *writer* receives a signal, it executes a *signalHandler* that determines the thread is non-restartable, and immediately returns. The *reclaimer* then goes on to reclaim its *limboBag* (line 17), except for any reserved records contained therein, independently from the actions of the *writer*.

This is safe because a *writer*, before entering into the  $\Phi_{write}$ , *reserves* all of the shared records it might access in its  $\Phi_{write}$  (line 11). Thus, (2) the *writer* guarantees to the *reclaimer* that, although it will not restart its data structure operation, it will only access *reserved* records. The (3) *reclaimer*, in turn, guarantees it will scan all announcements after signaling and before reclaiming the contents of its *limboBag*, and will consequently avoid reclaiming any records that will be accessed by the *writer* in its  $\Phi_{write}$ .

This three-step handshake formed by (1), (2) and (3) avoids scenarios where a *writer* might access a freed record. Crucially, all *writers* atomically ensure that their reserved records are visible to the *reclaimer* at the moment they become nonrestartable. In turn, *reclaimers* scan reservations *after* sending neutralization signals (at which point any thread that does not restart has already made its reservations visible).

#### 4.4 Revisiting the $\Phi_{write}$ restriction

In this section, we will trace an incorrect execution that could occur if a thread accesses any record that is not reserved *before* entering the  $\Phi_{write}$ .

Suppose a thread *T* is in a  $\Phi_{write}$ , and sleeps just before it accesses a shared record *rec*, which it has *not* reserved. Then, another thread *T'* sends a neutralization signal to *T* using signalAll. Next, *T'* scans the *reservations* array of the thread *T*. *T* did not reserve *rec* so *T'* will not find *rec* in *T's* reserved records (which violates the writers handshake, Section 4.3.2). Therefore, *T'* will assume that *rec* can be freed safely, and will do so. Finally, *T* wakes up and proceeds with its *unsafe* access of *rec*.

**Ensuring reclamation can occur.** The total number of records that can be reserved over all threads must be strictly smaller than the limboBag capacity, in order to ensure that threads *can* reclaim records whenever the limboBag is full. In practice, most data structures require few reservations. For example in our experiments, operations in the lazylist [31] required at most 2 reservations, and the Harris list [28], DGT [17], and (a,b)-trees [9] at most 3.

## 5 NBR+

Next, we explain a performance issue with NBR which motivated us to design an improved version called NBR+. Signals on linux trigger page-fault routines and a switch from user to kernel mode that incurs significant overhead. Therefore, it is desirable to send as few signals as possible (while maintaining high reclamation throughput). However, every time a thread reclaims records from its *limboBag*, NBR requires the thread to send signals to all other threads. This induces a relaxed grace period (RGP): A time interval [t, t'] during which each thread is neutralized due to a reclamation event triggered by some reclaimer thread. In NBR, every thread induces a RGP every time it tries to reclaim its limboBag. As a result, in order for all *n* threads to reclaim their *limboBags*, n(n-1) signals must be sent. The need to send  $O(n^2)$  signals to allow all *n* threads to reclaim memory can severely limit performance. Naturally, we would like to improve this.

Suppose, in NBR, at some time *t*, a thread sends n - 1signals to other threads so it can reclaim its *limboBag*. This causes all of the other n - 1 threads to discard any unreserved references to shared records. Meaning, at time t, the (unreserved) records in the limbo bags of all threads are safe to free. Therefore, if somehow we could propagate this information that a RGP has occurred due to some thread T, then all other threads could piggyback on T to partially or completely reclaim their own limbo bags without sending signals of their own. In other words, in the best case, all n participating threads could reclaim memory after detecting exactly one RGP, induced by sending a total of n - 1 signals. Overview of NBR+ The key insight in NBR+ is that when a reclaimer sends neutralization signals to all threads, all threads discard their pointers to unreserved records, and thus all threads can potentially reclaim some records in their *limboBags*. This suggests a design wherein each thread (1) passively detects a RGP by observing signals sent by another

**Algorithm 2** *NBR*+: Only variables that differ from *NBR* are shown here. *NBR*+ includes all variables and procedures of Algorithm 1. The retire operation is different in *NBR*+.

#### thread local variable:

1: int scanTS[N]; ▶ N = #threads. P = set of processes

- 2: bool firstLoWmEntryFlag = true;
- 3: record\* bookmarkTail;

#### shared variable:

4: atomic<int> announceTS[N];

#### 5: procedure retire(rec)

6:	if isAtHiWm() then	
7:	<pre>FAA(&amp;announceTS[tid],1);</pre>	▶ <i>RGP</i> begin
8:	signalAll()	
9:	<pre>FAA(&amp;announceTS[tid],1);</pre>	$\triangleright$ <i>RGP</i> end
10:	<pre>reclaimFreeable(tail);</pre>	
11:	<pre>cleanUp();</pre>	
12:	else if isAtLoWm() then	
13:	<pre>if firstLoWmEntryFlag then</pre>	
14:	<pre>bookmarkedTail = tail;</pre>	
15:	<pre>scanTS[tid] = scanAnnot</pre>	unceTS()
16:	end if	
17:	<b>for</b> each otid $\in$ P <b>do</b> $\triangleright$ otid: oth	er thread's id in P.
18:	<b>if</b> announceTS[otid]≥scanTS[	tid][otid]+2 then
19:	reclaimFreeable(book	kmarkTail);
20:	<pre>cleanUp();</pre>	
21:	break;	
22:	end if	
23:	end for	
24:	end if	
25:	limboBag[tid].append(rec);	
26:	end procedure	
27:	procedure CLEANUP	
28:	<pre>firstLoWmEntryFlag = 1;</pre>	
29:	end procedure	

thread, and (2) determines which records in its *limboBag* were unlinked *before* the RGP (i.e., are safe to reclaim).

### 5.1 Implementation of NBR+

We explain the design of *NBR+* by building our exposition around three main design challenges.

- (C1) When should a thread start tracking other threads' signals to detect a *RGP*?
- (C2) How can a thread *recognize* that a *RGP* has occurred?
- (C3) Once a thread recognises that a *RGP* has occurred how should it determine which records in its *limboBag* are safe to reclaim?

As a solution to **(C1)**, each thread in *NBR+*, in addition to watching the *limboBag* size to determine when it becomes *too large* (triggering neutralization), also determines when the *limboBag* size crosses a predetermined threshold called the *LoWatermark* (e.g., one half full or one quarter full). If

a thread's *limboBag* is full, we say that thread is at the *Hi-Watermark*. If a thread's *limboBag* keeps growing without reclamation it will first cross the *LoWatermark* and then hit the *HiWatermark*. As shown in Algorithm 2, a thread determines whether it has passed the *HiWatermark* or *LoWatermark* using procedures isAtHiWm (line 6) and isAtLoWm (line 12). Once a thread has passed the *LoWatermark*, it begins recording and analyzing information about signals sent by other threads to detect RGPs.

To tackle **(C2)**, a *reclaimer* at the *LoWatermark* (who wants to detect a *RGP*) must perform a sort of handshake with another *reclaimer* at the *HiWatermark* (who triggers a *RGP*). *NBR+* implements this handshake using per-thread single-writer multi-reader timestamps (similar to vector clocks).

Whenever a *reclaimer* hits the *HiWatermark*, it first increments its timestamp (to an odd value) to indicate that it is *currently broadcasting signals* (line 7). This denotes the *beginning* of a *RGP*. It then sends signals to all threads, and increments its timestamp again (to an even value) to indicate that it has *finished broadcasting signals* (line 9). This denotes the *end* of the *RGP*.

Whenever a reclaimer T passes the LoWatermark, it collects and saves the current timestamps of all threads (line 15), as well as the current *tail* pointer of its *limboBag* (line 14), so it can remember precisely which records it had unlinked before it reached its LoWatermark. T then periodically collects the timestamps of all threads, comparing the new values it sees to the original values it saw when it passed the LoWatermark (line 17 - line 23). (To obtain high performance, we amortize the overhead of scanning announceTS over many retire operations.) It continues to do this until it either detects a RGP or hits the HiWatermark itself (and sends signals to induce its own RGP). Observe that, after T hits its LoWatermark, if the timestamp of any thread changes from one even number to another even number, then that thread has both begun and finished sending signals to all threads since T hit the LoWatermark. Thus, T can identify that a RGP has occurred since it hit its LoWatermark, solving (C2).

Finally, to tackle **(C3)**, observe that T saves the last record (*tail* of its limboBag) it had retired before entering the *LoWa-termark* at line 14. If T successfully observes a *RGP* as explained in the solution to (C2), then all threads would either have discarded or reserved all their private references to the records in *T*'s limboBag up to the saved *bookmarkTail*. Thus, *T* can invoke reclaimFreeable to free all unreserved records up to the *bookmarkTail* (line 19). solving (C3).

cleanUp() (line 27) method is used to set *firstLoWmEntryFlag* after a thread reclaims either at *LoWatermark* (line 20) or at *HiWatermark* (line 11).

A thread that has not reached the *LoWatermark* or the *HiWatermark* simply continues to append any retired records to its *limboBag* (line 25).

At first it may appear that a thread *T* can reclaim its *lim*boBag as soon as it receives a neutralizing signal from a *reclaimer* thread *T'*. However, the receipt of a single signal is not enough for *T* to safely reclaim memory. To safely reclaim the set *R* of records in its *limboBag* up to its *bookmarkTail*, *T* needs to know that *all* threads have been neutralized *since T retired the records in R*. Otherwise, some other thread may still have a pointer to a record in *R*.

Let us discuss an example of what can go wrong if a thread reclaims its *limboBag* after it receives a single signal. Consider a system with three threads *T*1, *T*2 and *T*3. Suppose *T*1 is at its *HiWatermark*, T2 is at its *LoWatermark* and T3 holds a private reference to a record *rec* that is in *T2*'s *limboBag*. *T*1, being at its *HiWatermark*, begins neutralizing all threads one by one. First, it sends neutralizing signal to T2 (starting a *RGP*). *T*2, upon receiving the signal, reclaims its *limboBag* including rec. Note, that T1 hasn't neutralized T3 yet, meaning a *RGP* has not yet occurred. Now, if *T*3 accesses *rec*, a use-after-free error would occur. To prevent this, T2 should not reclaim the contents of its *limboBag* unless T1 completes the *RGP* by neutralizing *T*3 (preventing *T*3 from doing this unsafe access). The crucial point is that T2 must detect the start and end of a RGP to know that it can safely reclaim records in its *limboBag*.

### 5.2 Applicability

31 *NBR* (+)<sup>e</sup> naturally applies to many concurrent data struc-32 tures that have synchronization-free searches followed by 33 34 update(s) because in such data structures searches and updates correspond to the  $\Phi_{read}$  and the  $\Phi_{write}$  of NBR, respec-35 tively (as shown in Figure 1). Thus, to apply NBR one just 36 37 needs to invoke  $BEGIN\Phi_{read}$  before the start of the search 38 and  $END\Phi_{read}$  before starting the update(s). For example, in 39 the lazy-list of Heller et al. [31], the  $\Phi_{read}$  of an operation would begin with the start of the search for target records 41 and the  $\Phi_{write}$  would consist of the locking and validation 42 43 of target records followed by any modifications to them. 44

Certain other lock-free data structures exhibit a pattern where searches, in an operation, perform *auxiliary update(s)* followed by intended update(s). Such an operation has a *sequence* of *read-write phases*. For example, in Harris's lock-free list [28], while searching the list towards a target location, a thread may modify the list by unlinking any *marked* (logically deleted) records it encounters. Then, once it arrives at the target location, it performs the operation's intended modification.

Since *NBR* is designed for a single  $\Phi_{read}$  and  $\Phi_{write}$ , applying it carelessly to such a data structure could break the requirement that, after entering a  $\Phi_{write}$ , no new records are discovered. (This would be unsafe, because it breaks the *writers* handshake.) For instance, in such a data structure, if we enter a  $\Phi_{write}$  to perform an *auxiliary update*, *NBR* would be stuck in the  $\Phi_{write}$ , unable to obtain new pointers (that have not yet been reserved) to continue its traversal.

**Algorithm 3** Integration of *NBR* with Harris list[28] with multiple read/write phases  $(\Phi_{read}\Phi_{write})^+$ .

```
bool insert(key) {
  Node *right_node,
                     *left_node;
  do{
      right_node = search (key, &left_node);
      if((right_node!=tail) && (right_node.key==key))
        return false;
      Node *new_node = new Node(key);
      new_node.next = right_node;
      if (CAS(&(left_node.next), right_node, new_node))
        return true:
 }while (true)
}
Node* search(key, Node** left_node) {
  Node *left_node_next, *right_node;
  search_again:
  do {
      begin \Phi_{read}();
      Node *t = head;
      Node *t_next = head.next;
      do {
          if(!is_marked_reference(t_next)){
            (*left_node) = t;
            left_node_next = t_next;
            = get unmarked reference(t next);
          if (t == tail) break;
          t_next = t.next;
      }while(is_marked_reference(t_next) or (t.key<</pre>
           search_key));
      right_node = t;
      end\Phi_{read}(left_node, right_node);
      if (left_node_next == right_node)
        if ((right_node != tail) && is_marked_reference(
             right_node.next))
          goto search_again;
        else
          return right_node;
      if (CAS(&(left_node.next), left_node_next,
           right_node))
        if ((right_node != tail) && is_marked_reference(
             right_node.next))
          goto search_again;
        else
          return right_node;
 } while(true);
}
```

That said, *NBR can* be applied in some data structures that would require *multiple read/write* phases, provided that each consecutive pair of read and write phases obey the requirements set out in Section 4.1.

**Example:** Harris list. Algorithm 3 shows how *NBR* can be used with the Harris list [28], *despite* the fact that this list has auxiliary updates. We hope the reader can follow our exposition on the Harris list, and infer how *NBR* could be applied to more sophisticated data structures with similar design patterns (such as Brown's ABTree [9], which appears in our experiments).

To understand how *NBR* behaves when applied to the Harris list, suppose the initial list configuration is  $L: 1_f \implies 2_f \implies 3_t \implies 4_f \implies 6_f \implies 10_f$ , where each node is represented as  $key_{marked}$  (where *marked* is [*t*]rue or [*f*]alse). Now, suppose a thread *T* performs *Ins:insert(9)*, starting with an invocation of search(). This invocation of search()

10

11

12

14

16

17

18

20 21

23

24

25

26

27

28

29

30

<sup>&</sup>lt;sup>e</sup>We will simply write *NBR* in this section with the understanding that the entire discussion applies identically to *NBR+*.

starts a  $\Phi_{read}$  (line 18) and begins traversing *L*. Starting from  $\langle pred, curr \rangle = \langle 1_f, 2_f \rangle$  the thread observes  $\langle pred, curr \rangle = \langle 2_f, 3_t \rangle$ , where  $curr = 3_t$  is marked. To remove marked node  $3_t$  (an auxiliary *helping* update), *T* enters a  $\Phi_{write}$  (line 31) and changes the next pointer of  $2_f$  to  $4_f$ , yielding the list configuration:  $1_f \implies 2_f \implies 4_f \implies 6_f \implies 10_f$ . Moving forward, *T*'s search() will enter a *second*  $\Phi_{read}$  (line 18), and traverse the list again, *starting from the root*. As *T* now obtain pointers to *new nodes* (which would be impossible with only a single  $\Phi_{read}$  and  $\Phi_{write}$ ), we must argue that it doesn't access any freed nodes. However, this is straightforward, since it is again traversing *from the root* discarding any references from previous *read-write* phases. (From the perspective of SMR, it is as if *T* has simply started a new data structure operation.)

Now, suppose *T* is *neutralized* by a concurrent *reclaimer* while it is in this second  $\Phi_{read}$ . Upon receipt of a neutralization signal, *T* will jump back to the beginning of its *second*  $\Phi_{read}$ , and restart its search, once again, from the *root*. Note that neutralizing does not affect the lock-free progress guarantee, since a thread sends neutralization signals only after performing many successful deletion operations. Suppose *T* eventually performs a  $\Phi_{read}$  that reaches the nodes  $\langle pred, curr \rangle = \langle 6_f, 10_f \rangle$  where it should perform its modification. *T* will then enter a final  $\Phi_{write}$  and insert  $9_f$  after returning (line 37) from the search(), yielding *L*:  $1_f \Longrightarrow 2_f \Longrightarrow 4_f \Longrightarrow 6_f \Longrightarrow 9_f \Longrightarrow 10_f$ .

**Limitation: restarting from the root.** In order for *NBR* to be safe, it is *crucial that Ins forgets all pointers and restarts from the root every time it begins a new*  $\Phi_{read}$ . Intuitively, this is because each new read phase is effectively a new data structure operation—all pointers are forgotten when the new  $\Phi_{read}$  begins. If it attempts to continue searching from somewhere in the middle of the list, perhaps by resuming its search from a shared node *R* that was reserved by the previous  $\Phi_{write}$ , then *Ins* could easily dereference a freed node. To see why, note that, although *R* cannot be freed (since it is reserved), the nodes that it points to are *not* necessarily reserved, and so they could be freed. Thus, as soon as *Ins* follows any pointer starting from *R*, it could access a freed node and crash.

**Compatible data structures.** There are numerous concurrent data structures in the literature with multiple *read-write phases* that *do* restart from the *root* after any *auxiliary updates*, and hence are natural candidates for pairing with *NBR*. For example, Harris' list [28], Brown's lock-free AB-Tree, chromatic tree and AVL tree (B17) [9], the lock-free binary search tree of Natarajan et al. [36], and many more [23, 30, 33, 43]. Among these, we used the lazylist and the ABTree in our experiments. Due to space constraints, the Harris list appears in longer version of the paper[44].

**Semi-compatible data structures.** The need to restart from the root at the start of each  $\Phi_{read}$  suggests that *NBR* cannot be used with the data structures like the Harris-Michael list

[35], and some search trees [8, 12, 20, 22, 40], wherein the searches ( $\Phi_{read}$ ) after each *auxiliary update* ( $\Phi_{write}$ ) do not start from the *root*. However, we could potentially use *NBR* with such data structures if we were to modify the operations so they restart *from the root* after any auxiliary updates. Depending on the data structure, this might break the progress guarantee (for example changing a wait-free algorithm into a lock-free one, or necessitating a new amortized complexity analysis [22]), or it might simply add overhead.

For some data structures, the overhead of restarting from the *root* may be quite low in practice, and forcing operations to restart from the root may be a reasonable solution. (The cost of restarting from the root is studied in our experiments.) For example, in Harris-Michael list, in high contention scenarios where *k* threads all contend on an auxiliary CAS to unlink the same marked node, all threads except for the one that succeeds this CAS would already restart from the root [35] anyway! If we modify this list so threads always restart from the root, in this high contention scenario, *k* threads must restart instead of k - 1. Incidentally, by doing this, we essentially obtain the Harris list [28], in which all threads contending on the auxiliary CAS already restart from the *root*. (In low contention scenarios the way we restart should not affect performance significantly.)

Furthermore, in search trees, assuming a uniform distribution of accesses, threads tend to spread out in the tree, so average contention is quite low. This suggests that, when a thread encounters contention, the performance difference between restarting from the root and continuing a traversal from an ancestor will be small in many workloads.

**Incompatible data structures.** We are aware of a few data structures that are incompatible with (or would require extensive code changes to work with) *NBR*. Two concurrent implementations of a relaxed-balance AVL tree appear in [8, 20]. In each of these implementations, after a key is inserted, rotations must be performed to rebalance the tree. These rotations are performed starting at the bottom of the tree, possibly continuing all the way to the root (traversing upwards using parent pointers). In the process of performing these rotations, a thread may encounter many new nodes that were *not* traversed as part of the initial search in the insert operation. In order to use *NBR* with these algorithms, one would need to rewrite the implementations to perform a new search from the root *for each rotation*.

A recent lock-free interpolation search tree [12] also appears to be incompatible. For example, in this algorithm, entire subtrees are periodically rebuilt to maintain balance, and during this process, threads *mark* all nodes in the subtree, one by one, alternating between steps that *mark* a node and *discover* a new node (without restarting from the root in between). It is not clear how one could transform this algorithm into the form required by *NBR*. (Note, however, that neither DEBRA+ nor HPs can be used with this data

Source	Data structure	Sync. type	NBR+	EBR	DEBRA+	HP/TS/IBR/HE/WFE/HY/QSense
LL05[31]	linked list	opt. locks	Yes	Yes	No	No (similar to [13])
HL01[28]	linked list	lock-free	Yes	Yes	*	Yes
HM04[35]	linked list	lock-free	No	Yes	*	Yes
DVY14a[20]	partially external BST	locks	**	Yes	No	No [13]
EFRB10[23]	external BST	lock-free	Yes	Yes	*	No [13]
NM14[36]	external BST	lock-free	Yes	Yes	*	No [13]
EFRB14[22]	external BST	lock-free	No	Yes	*	No [13]
DGT15[17]	external BST	ticket locks	Yes	Yes	No	No (no marks, cannot validate HP)
HJ12[33]	internal BST	lock-free	Yes	Yes	*	No (similar to [13])
RM15[41]	internal BST	lock-free	No	Yes	No	No (similar to [13])
BCCO10[8]	partially external AVL	opt. locks	No	Yes	No	Yes
DVY14b[20]	partially external AVL	locks	No	Yes	No	No [13]
HL17[30]	external relaxed AVL tree	lock-free	Yes	Yes	Yes	No (similar to [13])
B17b[9]	external AVL	lock-free	Yes	Yes	Yes	No [13]
S13[43]	patricia trie	lock-free	Yes	Yes	*	No [13]
BER14[11]	external chromatic tree	lock-free	Yes	Yes	Yes	No [13]
B17a[9]	external (a,b)-tree	lock-free	Yes	Yes	Yes	No [13]
BPA20[12]	external interpolation tree	lock-free	No	Yes	No	No (similar to [13])

**Table 1. Applicability of SMR algorithms.** Due to space constraints a detailed explanation of the contents of this table is relegated to [44]. \*It appears likely that DEBRA+ is compatible, but one must design non-trivial data structure specific recovery code. \*\*This is likely possible if code is restructured to reserve all relevant nodes before acquiring any locks.

structure, either. We are not aware of *any* SMR algorithm with bounded garbage that is compatible with this tree.)

**Comparing with other SMR algorithms.** *NBR* can be used with many data structures that other SMR algorithms like DEBRA+ and HP (and variants of HPs, including HE, IBR, WFE, ThreadScan, HY and QSense) are incompatible with [17, 23, 31, 33, 36]. There are also some data structures that are compatible with other SMR algorithms but not *NBR* [8, 35]. See Table 1 for an overview. Due to lack of space, a detailed analysis of the table's contents, and an exposition of how *NBR* can be applied to these data structures, is relegated to the full version of this paper [44].

#### 5.3 Ease of use

Figure 2 compares the difficulty of using *HP*, *NBR* and *DEBRA* in the insert operation of the lazy list of Heller et al. [31]. As Figure 2c demonstrates, *HP* is cumbersome to use because it requires a programmer to protect every record by *announcing* hazard pointers, using a store/load fence or *xchg* instruction to ensure that each announcement is visible in a timely manner by other threads, validating that the announced record is still safe before dereferencing it, and restarting if validation fails. Programmers also need to unprotect records that they will no longer dereference, further increasing the need for intrusive code changes.

On the contrary, applying *NBR* to a data structure operation is, intuitively, similar to performing two-phased locking, in the sense that the primary difficulty revolves around identifying where the  $\Phi_{write}$  should begin, and which records it will access. The programmer just needs to invoke begin $\Phi_{read}$ before the operation accesses its first shared record, in this example, at the start of the traversal for target records. Then s/he must invoke end $\Phi_{read}$  before modifying any shared records. In this example, the  $\Phi_{write}$  begins just before the lock acquisition on pred. If there are no modifications to be performed in an operation, for example, in the contains operation of the lazy-list, then the programmer can simply invoke end $\Phi_{read}$  before returning from the operation.

*DEBRA* is simplest as it requires programmers to invoke just two functions corresponding to the start and the end of a data structure operation (Figure 2a).

In terms of programmer effort, *NBR* finds a middle ground between DEBRA and HPs. Although *NBR* is slightly more involved than DEBRA, we believe that the benefits due to *NBR*'s bounded garbage property and better performance outweigh the extra effort of identifying which shared records will be modified by the  $\Phi_{write}$  and where in the code to invoke end $\Phi_{read}$ .

Just to give readers a quantitative view of the amount of programming effort needed to use HP and NBR we measured number of extra reclamation related lines of code needed to be written in our implementation of insert(), delete() and contains() methods for the lazylist and DGT. We observed that *NBR* required only 10 extra lines of code in comparison to 30 extra lines of code needed to use HP.

As mentioned earlier, in *NBR*, before a thread enters a  $\Phi_{write}$ , it must reserve all the records that will be accessed in the  $\Phi_{write}$ . In some data structures it might not be possible to determine *precisely* which records *will* be accessed in the  $\Phi_{write}$ . For example, in a tree, an operation may decide *during* the write phase whether to modify the left or right child pointer. To apply *NBR* in such a tree, one can simply



Figure 2. Complexity of using DEBRA, NBR and HP on a lazy list. DEBRA << NBR << HP.

reserve *both* pointers. (Nevertheless there may be some data structures where it is infeasible to reserve *all* of the records that might be accessed in a  $\Phi_{write}$ .)

#### 5.4 Correctness

*NBR* and *NBR*+ are both safe and have bounded garbage. Proofs can be found in [44].

## 6 Experimental Evaluation

**Setup:** We used a quad-socket Intel Xeon Platinum 8160 machine running at 2.1GHz with 192 hardware threads and 384 GiB memory having shared L3 cache (33.79 MiB) on Ubuntu 18.04 with GCC/G++ version 7.4.0.

All algorithms used in the experiments were implemented in the Setbench [12] benchmark compiled with -03 optimization, and used *jemalloc* as the memory allocator [24]. We perform four kinds of experiments:

- (E1): Studies throughput over different thread counts and workloads to understand *NBR* (+)'s performance and scalability [P1].
- (E2): Studies peak memory usage of *NBR+* to show the advantage of bounded garbage [P2].
- (E3): Studies the impact of contention on performance [P4].
- (E4): Studies the impact of modifying a data structure to restart from the root to make it work with *NBR+*.

For (E1) we picked the lazy-list [31] and DGT [17] as representative of lists and trees, to evaluate *NBR+* against QSBR, RCU, DEBRA, and the 2geibr variant of interval based reclamation (IBR) [46], hazard pointers (HP) and a leaky implementation (none). (We adapted QSBR, RCU and 2geibr (IBR) from the IBR benchmark, integrating them into Setbench to ensure a fair comparison.) For (E2) we compared peak memory usage for each of the aforementioned reclamation algorithms using DGT. For (E3) we compared *NBR* with DEBRA and the leaky implementation (none) using the ABTree data structure [9, chapter 8] for very large and small data structure size. Finally, for (E4), we modified the Harris-Michael list (HMList) such that every  $\Phi_{read}$  restarts from the root. This allowed us to use *NBR+* to reclaim memory for this list. To understand the impact of these restarts, we also used DEBRA to reclaim memory for this modified HMList (labeled *DEBRA-restarts* in our graphs), as well as the original HMList (DEBRA-norestarts).

Reported results are obtained by averaging data over 3 timed trials, each lasting 5 seconds, for each thread count in {24, 48, 72, 96, 120, 144, 168, 192, 216, 240, 252}, and each data structure. We used a key range size of 2 M and 20 K for trees and lists, respectively. Each execution starts by prefilling the data structure to half of the key range size, i.e., 1 M for trees and 10 K for lists.

For each of (E1), (E2), and (E3) we subject *NBR+* to exhaustive evaluation by running it on three workload profiles, (1) *update-intensive* 50% inserts and 50% deletes, (2) *balanced* 25% inserts, 25% deletes and 50% searches, and (3) *searchintensive* 90% searches, 5% inserts and 5% deletes. (We also run with oversubscription—i.e., more threads than logical processors, to establish P4.)

**Discussion** (E1) results for the DGT tree and lazy linked list appear in Figure 3. *NBR+* matches or outperforms its competitors for nearly all data points. In the tree, it surpasses the next best algorithm, DEBRA, by up to 38% and 12% (Figure 3a, update-intensive and balanced workloads, resp.) and is comparable to DEBRA in search-intensive workloads where it

#### Ajay Singh, Trevor Brown, and Ali Mashtizadeh







(b) Lazy linked-list. Left: 50i-50d. Middle: 25i-25d. Right: 5i-5d. Key range size:20K.

Figure 3. Evaluation of throughput. Y axis: throughput in million operations per second. X axis: #threads.

outperforms the next best algorithm by up to 10% (Figure 3a search-intensive workload).

In these graphs, DEBRA performs better than NBR+ for low thread counts, but NBR+ outperforms it after 96 threads in update intensive workloads (Figure 3a, leftmost plot), and after 120 threads in the balanced workload (Figure 3a, center). The two algorithms are comparable in the search-intensive workload (Figure 3a, rightmost plot). The poor performance of DEBRA at higher thread counts could be attributed to the infrequent advancement of epochs by slow threads, which leads to halting of regular reclamation of limbo bags. We call this the *delayed thread vulnerability*. As a result, the limbo bags of all threads keep on growing until the slow thread finally catches up and announces the required epoch.

The delayed thread vulnerability leads to the accumulation of a large number of retired records waiting to be reclaimed. Once the current epoch is announced by the slow thread, all threads reclaim their large limbo bags, which leads to a *reclamation burst*. This harms the overall throughput, as reclamation bursts can bottleneck the underlying allocator (jemalloc in our experiments) by increasing contention and triggering slow/fallback code paths. The probability of threads getting delayed increases with high thread counts and update-intensive workloads.

Furthermore, one may notice that the thread count where *NBR+* overtakes DEBRA is different in the *update intensive* and *search-intensive* workloads. This could be attributed to the fact that the overhead of burst reclamation sets in at lower thread counts for *update-intensive* workloads than in workloads with infrequent updates.

HP outperforms the other EBR variants in update-intensive workloads (Figure 3a, leftmost plot) but they appear to be slow in the search-intensive workload (Figure 3a, rightmost plot). This can be attributed to the fact that overhead due to the delayed thread vulnerability dominates the overhead of per-record fencing in HP for update-intensive workloads. Whereas, for search-intensive workloads the overhead due to delayed threads is lesser than that of per-record fences in HPs. Also, as one can observe in Figure 3b, in the lazylist *NBR+* is comparable to *RCU*, *QSBR*, and *DEBRA*, and performs better than *HP* (by 2x) and *IBR* (by more than 50%) across all workloads when oversubscribed [P1, P4].<sup>f</sup>

In E2 (Figure 4c and 4d), we validate the bounded garbage property of NBR+ [P2] by measuring the peak memory usage of all reclamation algorithms when a thread is stalled (Figure 4c) and when no thread is stalled (Figure 4d). Each trial is run for 25 seconds. During the entire length of the experiment (25 seconds) one thread is made to begin a data structure operation and then sleep. As expected, since *DE-BRA* and *RCU* do not have bounded garbage, they exhibit an increase in peak memory usage in presence of a stalled thread (Figure 4c) while *NBR+*, *HP*, and *IBR* variants maintain approximately the same peak memory usage.

In E3 (Figure 4a) we evaluate throughput of *NBR* with a data structure in which operations already restart from the root whenever a new read phase would be started, namely the Brown ABTree[9, chapter 8]. We design our experiments

<sup>&</sup>lt;sup>f</sup>The poor performance of HPs is partly due to the cost of *mfence*, which could be replaced with the more efficient xchg (see Section 11.5.1 in https://www.amd.com/system/files/TechDocs/47414\_15h\_sw\_opt\_guide.pdf). We tried this optimization in separate experiments, and while it did increase throughput, HPs remained significantly slower than *NBR+*.



(a) Brown's ABTree. Left: 50i-50d. Key range size: 2M. Right: 50i-50d. Key range size:200.







(b) Lock-free Harris-Michael list. Left: 50i-50d. Key range size:20K. Right: 50i-50d. Key range size:200. The debra-norestarts use HMlist[35] whereas *NBR+*, debra-restarts and none use modified HMList that restarts from the root.

(d) NO stalled threads.

**Figure 4.** Fig. (a), (b): Evaluation of throughput for data structures with multiple read-write phases. Y axis: millions of ops/sec. Fig. (c), (d): Y axis: Max Resident Memory (MB), Workload: 50i-50d. DGT-tree. Key range size: 2M.

to explore two disparate usage scenarios. First, we want to study reasonably large data structure (key range size 2 M) wherein we hypothesize restarts would be inexpensive due to low contention (leftmost in Figure 4). Second, we want to study extremely small data structures (key range size 200) where restarting from the head node will occur as frequently as possible due to high contention (rightmost in Figure 4).

In the ABTree, *NBR+* is faster than the other SMR algorithms in the low contention scenario (especially at high thread counts). Moreover, in the high contention scenario, *NBR+* is comparable to DEBRA, suggesting that *NBR+* introduces relatively little overhead in practice due to restarting from the root in  $\Phi_{read}$  (Figure 4a).

In E4 (Figure 4b), for low contention, *NBR+* (with forced restarts from the root) is faster than both DEBRA-restarts and DEBRA-norestarts. Shockingly, DEBRA-restarts is actually slightly *faster* than DEBRA-norestarts. The only code difference between these two implementations is an extra restart from the root in one case of DEBRA-restarts. It appears that forced restarts actually have a backoff-like (contention managing) effect. We found that adding restarts lowered L3 cache misses slightly, which is what we would expect from contention management. Under high contention, *NBR+* is comparable to DEBRA-restarts and DEBRA-norestarts. This suggests that the cost of added restarts should be reasonable in practice.

**Code:** The latest version of our source code can be found on gitlab (https://gitlab.com/aajayssingh/nbr\_setbench).

## 7 Conclusions

In this paper, we presented *NBR*, a *safe memory reclamation* algorithm that is a hybrid between EBR and a limited form of HPBR, and which uses POSIX signals to bound unreclaimed garbage. *NBR* is simpler to use than the most similar hybrid algorithm, DEBRA+, while supporting a large class of data structures, some of which are not supported by DEBRA+ (nor other popular SMR algorithms). We also developed an optimized version of *NBR* called *NBR+* that achieves similar throughput with fewer signals by passively observing signals being sent in the system to optimistically detect relaxed grace periods. Our experiments demonstrate that *NBR+* meets or exceeds the performance of the state of the art in SMR algorithms in typical benchmark conditions, while minimizing performance degradation in oversubscribed workloads.

## Acknowledgments

This work was supported by the Natural Sciences and Engineering Research Council of Canada (NSERC) under the grants: CRDPJ 539431-19, DGECR-2019-00048, RGPIN-2019-04227 and RGPIN-04512-2018. The John R. Evans Leaders Fund and the Ontario Research Fund (CFI):38512, Waterloo Huawei Joint Innovation Lab project "Scalable Infrastructure for Next Generation Data Management Systems", and the University of Waterloo. We would also like to thank the reviewers for their insightful comments. PPoPP '21, February 27-March 3, 2021, Virtual Event, Republic of Korea

#### Ajay Singh, Trevor Brown, and Ali Mashtizadeh

## References

- Yehuda Afek, Haim Kaplan, Boris Korenfeld, Adam Morrison, and Robert E Tarjan. 2014. The CB tree: a practical concurrent selfadjusting search tree. *Distributed computing* 27, 6 (2014), 393–417.
- [2] Dan Alistarh, Patrick Eugster, Maurice Herlihy, Alexander Matveev, and Nir Shavit. 2014. Stacktrack: An automated transactional approach to concurrent memory reclamation. In Proceedings of the Ninth European Conference on Computer Systems. 1–14.
- [3] Dan Alistarh, William Leiserson, Alexander Matveev, and Nir Shavit. 2017. Forkscan: Conservative memory reclamation for modern operating systems. In *Proceedings of the Twelfth European Conference on Computer Systems*. 483–498.
- [4] Dan Alistarh, William Leiserson, Alexander Matveev, and Nir Shavit. 2018. Threadscan: Automatic and scalable memory reclamation. ACM Transactions on Parallel Computing (TOPC) 4, 4 (2018), 1–18.
- [5] Oana Balmau, Rachid Guerraoui, Maurice Herlihy, and Igor Zablotchi. 2016. Fast and robust memory reclamation for concurrent data structures. In *Proceedings of the 28th ACM Symposium on Parallelism in Algorithms and Architectures*. 349–359.
- [6] Guy E Blelloch and Yuanhao Wei. 2020. Concurrent Reference Counting and Resource Management in Wait-free Constant Time. arXiv preprint arXiv:2002.07053 (2020).
- [7] Anastasia Braginsky, Alex Kogan, and Erez Petrank. 2013. Drop the anchor: lightweight memory management for non-blocking data structures. In Proceedings of the twenty-fifth annual ACM symposium on Parallelism in algorithms and architectures. 33–42.
- [8] Nathan G Bronson, Jared Casper, Hassan Chafi, and Kunle Olukotun. 2010. A practical concurrent binary search tree. ACM Sigplan Notices 45, 5 (2010), 257–268.
- [9] Trevor Brown. 2017. Techniques for Constructing Efficient Lock-free Data Structures. arXiv preprint arXiv:1712.05406 (2017).
- [10] Trevor Brown, Faith Ellen, and Eric Ruppert. 2014. A General Technique for Non-blocking Trees. In *Proceedings of the 19th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP* '14). 329–342. Full version available from http://tbrown.pro.
- [11] Trevor Brown, Faith Ellen, and Eric Ruppert. 2014. A general technique for non-blocking trees. In Proceedings of the 19th ACM SIGPLAN symposium on Principles and practice of parallel programming. 329–342.
- [12] Trevor Brown, Aleksandar Prokopec, and Dan Alistarh. 2020. Nonblocking interpolation search trees with doubly-logarithmic running time. In Proceedings of the 25th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming. 276–291.
- [13] Trevor Alexander Brown. 2015. Reclaiming memory for lock-free data structures: There has to be a better way. In *Proceedings of the 2015* ACM Symposium on Principles of Distributed Computing. 261–270.
- [14] Nachshon Cohen. 2018. Every data structure deserves lock-free memory reclamation. *Proceedings of the ACM on Programming Languages* 2, OOPSLA (2018), 1–24.
- [15] Nachshon Cohen and Erez Petrank. 2015. Automatic memory reclamation for lock-free data structures. ACM SIGPLAN Notices 50, 10 (2015), 260–279.
- [16] Nachshon Cohen and Erez Petrank. 2015. Efficient memory management for lock-free data structures with optimistic access. In *Proceedings* of the 27th ACM symposium on Parallelism in Algorithms and Architectures. 254–263.
- [17] Tudor David, Rachid Guerraoui, and Vasileios Trigonakis. 2015. Asynchronized concurrency: The secret to scaling concurrent search data structures. ACM SIGARCH Computer Architecture News 43, 1 (2015), 631–644.
- [18] David L Detlefs, Paul A Martin, Mark Moir, and Guy L Steele Jr. 2002. Lock-free reference counting. *Distributed Computing* 15, 4 (2002), 255–271.

- [19] Dave Dice, Maurice Herlihy, and Alex Kogan. 2016. Fast non-intrusive memory reclamation for highly-concurrent data structures. In Proceedings of the 2016 ACM SIGPLAN International Symposium on Memory Management. 36–45.
- [20] Dana Drachsler, Martin Vechev, and Eran Yahav. 2014. Practical concurrent binary search trees via logical ordering. In Proceedings of the 19th ACM SIGPLAN symposium on Principles and practice of parallel programming. 343–356.
- [21] Aleksandar Dragojević, Maurice Herlihy, Yossi Lev, and Mark Moir. 2011. On the power of hardware transactional memory to simplify memory management. In *Proceedings of the 30th annual ACM SIGACT-SIGOPS symposium on Principles of distributed computing*. 99–108.
- [22] Faith Ellen, Panagiota Fatourou, Joanna Helga, and Eric Ruppert. 2014. The amortized complexity of non-blocking binary search trees. In Proceedings of the 2014 ACM symposium on Principles of distributed computing. 332–340.
- [23] Faith Ellen, Panagiota Fatourou, Eric Ruppert, and Franck van Breugel. 2010. Non-blocking binary search trees. In *Proceedings of the 29th ACM SIGACT-SIGOPS symposium on Principles of distributed computing*. 131– 140.
- [24] Jason Evans. 2006. A scalable concurrent malloc (3) implementation for FreeBSD. In *Proc. of the bsdcan conference, ottawa, canada.*
- [25] Panagiota Fatourou, Elias Papavasileiou, and Eric Ruppert. 2019. Persistent non-blocking binary search trees supporting wait-free range queries. In *The 31st ACM Symposium on Parallelism in Algorithms and Architectures*. 275–286.
- [26] Keir Fraser. 2004. Practical lock-freedom. Technical Report. University of Cambridge, Computer Laboratory.
- [27] Anders Gidenstam, Marina Papatriantafilou, Håkan Sundell, and Philippas Tsigas. 2008. Efficient and reliable lock-free memory reclamation based on reference counting. *IEEE Transactions on Parallel and Distributed Systems* 20, 8 (2008), 1173–1187.
- [28] Timothy L Harris. 2001. A pragmatic implementation of non-blocking linked-lists. In *International Symposium on Distributed Computing*. Springer, 300–314.
- [29] Thomas E Hart, Paul E McKenney, Angela Demke Brown, and Jonathan Walpole. 2007. Performance of memory reclamation for lockless synchronization. J. Parallel and Distrib. Comput. 67, 12 (2007), 1270–1285.
- [30] Meng He and Mengdu Li. 2017. Deletion without rebalancing in non-blocking binary search trees. In 20th International Conference on Principles of Distributed Systems (OPODIS 2016). Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik.
- [31] Steve Heller, Maurice Herlihy, Victor Luchangco, Mark Moir, William N Scherer, and Nir Shavit. 2005. A lazy concurrent list-based set algorithm. In *International Conference On Principles Of Distributed Systems*. Springer, 3–16.
- [32] Maurice Herlihy, Victor Luchangco, Paul Martin, and Mark Moir. 2005. Nonblocking memory management support for dynamic-sized data structures. ACM Transactions on Computer Systems (TOCS) 23, 2 (2005), 146–196.
- [33] Shane V Howley and Jeremy Jones. 2012. A non-blocking internal binary search tree. In Proceedings of the twenty-fourth annual ACM symposium on Parallelism in algorithms and architectures. 161–171.
- [34] Paul E McKenney and John D Slingwine. 1998. Read-copy update: Using execution history to solve concurrency problems. In *Parallel and Distributed Computing and Systems*, Vol. 509518.
- [35] Maged M Michael. 2004. Hazard pointers: Safe memory reclamation for lock-free objects. *IEEE Transactions on Parallel and Distributed Systems* 15, 6 (2004), 491–504.
- [36] Aravind Natarajan and Neeraj Mittal. 2014. Fast concurrent lockfree binary search trees. In Proceedings of the 19th ACM SIGPLAN symposium on Principles and practice of parallel programming. 317– 328.

- [37] Ruslan Nikolaev and Binoy Ravindran. 2019. Hyaline: fast and transparent lock-free memory reclamation. In *Proceedings of the 2019 ACM* Symposium on Principles of Distributed Computing. 419–421.
- [38] Ruslan Nikolaev and Binoy Ravindran. 2020. Universal wait-free memory reclamation. In Proceedings of the 25th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming. 130–143.
- [39] Aleksandar Prokopec, Nathan Grasso Bronson, Phil Bagwell, and Martin Odersky. 2012. Concurrent tries with efficient non-blocking snapshots. In Proceedings of the 17th ACM SIGPLAN symposium on Principles and Practice of Parallel Programming. 151–160.
- [40] Arunmoezhi Ramachandran and Neeraj Mittal. 2015. CASTLE: fast concurrent internal binary search tree using edge-based locking. ACM SIGPLAN Notices 50, 8 (2015), 281–282.
- [41] Arunmoezhi Ramachandran and Neeraj Mittal. 2015. A fast lock-free internal binary search tree. In *Proceedings of the 2015 International Conference on Distributed Computing and Networking*. 1–10.
- [42] Pedro Ramalhete and Andreia Correia. 2017. Brief announcement: Hazard eras-non-blocking memory reclamation. In Proceedings of the 29th ACM Symposium on Parallelism in Algorithms and Architectures. 367–369.
- [43] Niloufar Shafiei. 2013. Non-blocking Patricia tries with replace operations. In 2013 IEEE 33rd International Conference on Distributed Computing Systems. IEEE, 216–225.
- [44] Ajay Singh, Trevor Brown, and Ali Mashtizadeh. 2020. NBR: Neutralization Based Reclamation. arXiv:2012.14542 [cs.DC]
- [45] Shahar Timnat and Erez Petrank. 2014. A practical wait-free simulation for lock-free data structures. ACM SIGPLAN Notices 49, 8 (2014), 357– 368.
- [46] Haosen Wen, Joseph Izraelevitz, Wentao Cai, H Alan Beadle, and Michael L Scott. 2018. Interval-based memory reclamation. ACM SIGPLAN Notices 53, 1 (2018), 1–13.

## 8 Artifact Description

This section provides a step by step guide to run our artifact (nbr\_setbench) in a docker container. Other ways to setup a machine to use our artifact are provided in the *readme* file at the URL: https://doi.org/10.5281/zenodo.4409185.

To better reproduce the results described in our paper we recommend to run the nbr\_setbench on a NUMA machine with atleast 2 NUMA nodes having a recent Linux distro (we used Ubuntu 18.04 or 20.04) with 188GB RAM and recent Docker installation (we used version 19.03.6, build 369ce74a3c).

**Steps to load and run the provided Docker image:** Sudo permission may be required to execute the following instructions.

1. Install the latest version of Docker on your system. We tested the artifact with the Docker version 19.03.6, build 369ce74a3c. Instructions to install Docker may be found at

https://docs.docker.com/engine/install/ubuntu.

\$ docker -v

- 2. Download the artifact from Zenodo at URL: https://doi.org/10.5281/zenodo.4409185.
- Extract the downloaded folder and move to nbr\_setbench/ directory using cd command.

4. Find docker image named *nbr\_docker.tar.gz* in nbr\_setbench/ directory. And load the downloaded docker image with the following command:

\$ sudo docker load -i nbr\_docker.tar.gz

5. Verify that image was loaded:

\$ sudo docker images

- 6. Start a docker container from the loaded image:
  - $\$  sudo docker run --name nbr -i -t  $\$
  - --privileged nbr\_setbench /bin/bash
- 7. Invoke *ls* to see several files/folders of the artifact: Dockerfile, README.md, common, ds, install.sh, lib, microbench, nbr\_experiments, tools.

**Steps to run the experiments:** To compile, run and see results of the experiment follow these steps:

**Input**: Inputs to the experiment can be configured in corresponding input files at:

/nbr\_setbench/nbr\_experiments/inputs/

**Output**: Generated figures can be found in directory: /nbr\_setbench/nbr\_experiments/plots/generated\_plots/

- 1. Assuming you are currently in *nbr\_setbench*, execute the following command:
  - \$ cd nbr\_experiments/
- 2. Run the following command to generate plots for throughput evaluation:
  - \$ ./run.sh
- 3. Run the following command to generate plots for memory usage evaluation:
  - \$ ./run\_memusage.sh

After the above scripts finish executing DO NOT exit the terminal as we would need to copy the generated figures on the host machine to be able to see them.

**Steps to visualize the plots:** Resultant figures could be found in *nbr experiments/plots/generated plots.* 

To visualize the generated figures on your host machine copy the plots from the docker container to your host system by following these steps:

- 1. Verify the name of the docker container. Use the following command which would give us the name of the loaded docker container under NAMES column which is 'nbr'.
  - \$ sudo docker container ls
- 2. Open a new terminal on the same machine. Move to any directory where you would want the generated plots to be copied (use cd). And execute the following command to copy the generated plots from the *nbr\_experiments/plots/generated\_plots* folder to your current directory.
  - \$ sudo docker cp nbr:/nbr\_setbench/ \
     nbr\_experiments/plots/generated\_plots/ .

PPoPP '21, February 27-March 3, 2021, Virtual Event, Republic of Korea

Now you can analyse the generated plots. Each plot follows a naming convention:

1. throughput-[data structure name]-[number of inserts]-[number of deletes].png. For example, a plot showing throughput of DGT with 50% inserts and 50% deletes is named as: throughput-guerraoui\_ext\_bst\_ticket-i50d50.png. 2. Similarly the plot for peak memory usage experiments follows a naming convention: mem\_usage-[data structure name]-[number of inserts]-[number of deletes].png. For example, a plot showing mem\_usage of DGT with 50% inserts and 50% deletes is named as: mem\_usage-guerraoui\_ext\_bst\_ticket-i50d50.png.