

# A Family of Fast and Memory Efficient Lock- and Wait-Free Reclamation

Supplemental Material (Anonymized)

## 1 Extended Correctness Discussion

Crystalline-L/-LW/-W are based on Hyaline-1S, for which correctness is discussed in [2]. We present memory bounds arguments which are critical for lock-free guarantees. We further discuss aspects that pertain to wait-free progress.

**Theorem 1.** *Crystalline-L and Crystalline-W are fully memory bounded.*<sup>1</sup>

*Proof.* Batch sizes are constrained by  $\text{MAX\_THREADS} \times \text{MAX\_IDX} + 1$  nodes. In the worst case, each batch is attached to every reservation. The total number of reservations is  $\text{MAX\_THREADS} \times \text{MAX\_IDX}$ . Consequently, the memory usage is bounded by  $(\text{MAX\_THREADS} \times \text{MAX\_IDX} + 1)^2$ . Since era increments are also amortized, the total cost is  $\text{RETIRE\_FREQ} \times (\text{MAX\_THREADS} \times \text{MAX\_IDX} + 1)^2$ . (Note that for Crystalline-W,  $\text{MAX\_IDX} + 2$  should be used rather than  $\text{MAX\_IDX}$  since the total number of indices is bigger.)

Note that this theoretical upper bound is worse than that of HE [3] or WFE [1]. However, in practice, batches do not accumulate this (worst-case) number of nodes and are retired much faster, often resulting in better practical efficiency. Regardless of that, this worst-case bound is still reasonable and finite.  $\square$

**Lemma 1.** *traverse() calls are wait-free bounded.*

*Proof.* The loop in traverse() is bounded by the length of the list of retired batches at a given reservation. try\_retire() attaches (Line 58, Figure 2) only those batches for which the minimum birth era overlaps with the reservation's era. (Note that every time the era is updated by update\_era(), the list is emptied.) Since the eras are periodically incremented in alloc\_node(), the number of such nodes, and consequently – batches, is finite.  $\square$

**Lemma 2.** *detach\_nodes()'s loop is bounded by MAX\_THREADS iterations.*

*Proof.* detach\_nodes() makes the tag- $\rightarrow$ tag+1 (odd) transition in the slow path. The CAS operation in Line 24, Figure 3, changes the era tag unless it was already changed by a concurrent thread. For the loop in Lines 26-31, Figure 3 to continue, the list tag should have not yet moved to the tag+1 state. Because the era tag moves to the tag+2 (even) state only after that transition (Line 43 or Lines 83-87, Figure 4), i.e., after detach\_nodes() in Line 38 or 81 (Figure 4), the era tag is still tag+1. Consequently, all contending threads in try\_retire() will find that the era tag is odd (Line 46, Figure 2) and will skip the corresponding reservation for retirement. (Note that skipping nodes is safe because this race window is handled later by Lines 46 and 89, Figure 4.) Only threads that are already in-progress will proceed and potentially contend because of unconditional list pointer updates in Line 58, Figure 2. The number of such threads is bounded by  $\text{MAX\_THREADS}$ , as subsequent try\_retire() calls will observe that the era tag is odd.  $\square$

**Lemma 3.** *The loop in Lines 92-99, Figure 4 is bounded by at most MAX\_THREADS iterations.*

---

<sup>1</sup>This is not generally true for Hyaline-1S due to starving threads.

*Proof.* The proof is similar to that of Lemma 2. The only difference is that it makes the tag+1→tag+2 (even) transition in the slow path. Consequently, `try_retire()` will find that the list tag is odd (Line 44, Figure 2).  $\square$

**Lemma 4.** *The loop in Lines 13-37, Figure 4 is bounded by at most  $\text{MAX\_THREADS}$  iterations.*

*Proof.* In Line 11, Figure 4, a thread advertises that it needs help. The loop in Lines 13-37 can only fail to converge because of `global_era` updates. At most  $\text{MAX\_THREADS}$  already in-progress threads are executing `increment_era()` from `alloc_block()`, prior to Line 35, Figure 2, which updates `global_era`, but after Line 32, Figure 2, which detects what threads need helping. All these threads will execute Line 35. That will cause the loop in Lines 13-37 (Figure 4) fail and repeat. However, all newer `increment_era()` calls will only update `global_era` after `help_thread()` is complete.  $\square$

**Lemma 5.** *The loop in Lines 74-106, Figure 4 is bounded by at most  $\text{MAX\_THREADS}$  iterations.*

*Proof.* The same idea as in Lemma 4. We also need to make sure that the loop will not go beyond one slow path cycle. This is achieved by comparing the tag component in Line 106.  $\square$

**Theorem 2.** *`retire()` is wait-free bounded.*

*Proof.* `retire()` periodically calls `try_retire()`. The loop in Lines 55-67, Figure 2, is bounded by the number of nodes in a batch. ( $\text{MAX\_THREADS} \times \text{MAX\_IDX} + 1$  at most). Extra  $\text{RETIRE\_FREQ} - 1$  nodes can be retired since `retire()` calls `try_retire()` with the corresponding frequency. Regardless of the status of the CAS operations (Lines 61 and 64), the loop moves on to the next node. The `traverse()` call (Line 65) is bounded due to Lemma 1.  $\square$

**Theorem 3.** *`alloc_node()` is wait-free bounded.*

*Proof.* `alloc_node()` calls `increment_era()`. The latter includes a bounded loop with calls to `help_thread()`. Finally, `help_thread()` is bounded due to Lemmas 2, 3, 5.  $\square$

**Theorem 4.** *`protect()` is wait-free bounded.*

*Proof.* The fast path includes a finite number of iterations. It may call `update_era()`, which calls `traverse()`. `traverse()` is bounded due to Lemma 1. `slow_path()` contains a loop which is bounded due to Lemma 4. `slow_path()` also calls `detach_nodes()`, which is bounded due to Lemma 2. Finally, `slow_path()` can call `traverse()`, which is bounded due to Lemma 1.  $\square$

**Theorem 5.** *`clear()` is wait-free bounded.*

*Proof.* The method has a bounded loop which can call `traverse()`, which is bounded due to Lemma 1.  $\square$

## 2 Crystalline-W's Extended Discussion and Pseudocode

Crystalline-W's high-level changes (with respect to Crystalline-L) are shown in Figure 1. Note that reservation's list and era are now tagged. Tags are used in slow-path procedures only; fast-path procedures simply use the value component. Crystalline-W defines per-thread *state* (for each corresponding reservation) used in slow-path procedures and *slow\_counter* to identify if any thread needs helping. Those are somewhat similar to WFE's [1] corresponding slow-path variables. Finally, Crystalline-W defines the *parent* array to facilitate object handover, as discussed previously. Object handover is unique to Crystalline-W since it cannot simply scan the list of retired objects twice, as WFE, to avoid race conditions. Figure 1 also modifies `alloc_node()` to internally call `increment_era()` in lieu of doing FAA on the global era directly. Finally, `protect()` calls `slow_path()` if it fails to converge after  $\text{MAX\_TRIES}$ .

Figure 2 shows changes to `try_retire()` and `traverse()`. These methods use list tainting, as previously discussed. Also, unlike WFE [1], we use two slow-path tag transitions (odd and even). This is needed to make

```

1  template <typename type> struct Tag {
2      type V; // Value
3      uint64 T; // Tag, or Era for State::result
4  };
5
6  struct Reservation {
7      Tag<Node*> list; // Init: {.V = nullptr, .T = 0}
8      Tag<uint64> era; // Init: {.V = 0, .T = 0}
9  };
10
11 struct State {
12     Tag<void*> result; // Init: {.V = nullptr, .T = 0}
13     uint64 era; // Init: 0
14     Node* parent; // Init: nullptr
15     Node** obj; // Init: nullptr
16 };
17
18 Reservation rsrv[MAX_THREADS][MAX_IDX+2];
19 State state[MAX_THREADS][MAX_IDX];
20 Node* parents[MAX_THREADS]; // Init: (all) nullptr
21 int slow_counter = 0;

```

```

22 // Help other threads before incrementing the era
23 Node* alloc_node(int size) {
24     if (alloc_cnt++ % ALLOC_FREQ == 0) increment_era();
25     Node* node = malloc(size);
26     node->birth = global_era;
27     node->blink = nullptr; // Retired if != nullptr
28     return node;
29 }
30
31 // Use the fast-path-slow-path method
32 Node* protect(Node** obj, int index, Node* parent) {
33     int tries = MAX_TRIES;
34     uint64 prev_era = rsrv[TID][index].era.V;
35     while (--tries != 0) {
36         Node* ptr = *obj;
37         uint64 curr_era = global_era;
38         if (prev_era == curr_era) return ptr;
39         prev_era = update_era(curr_era, index);
40     }
41     return slow_path(obj, index, parent);
42 }

```

Figure 1: Crystalline-W (API function changes).

the number of iterations finite in some loops (Lemmas 2 and 3) by collaborating with `try_retire()` which will skip odd tags. Figure 2 also shows `increment_era()`'s implementation as used by `alloc_node()`.

Crystalline-W's slow-path and helper thread routines are demonstrated in Figure 4. These routines use several utility methods shown in Figure 3. The idea is similar to that of WFE [1], with one major difference: we use the *parent* array to keep parent references to facilitate object hand-overs, as previously discussed.

Utility methods in Figure 3 are needed to facilitate the slow path. `get_birth_era()` uses a trick to retrieve the birth era irrespective of whether the node is retired. WFE [1] always keeps the birth era. However, the Hyaline and Crystalline schemes recycle the birth era field after retirement so that they still use 3 words per each memory object. Since the birth era does not survive node retirements (except REFS which stores the minimum era for the batch), that presents a challenge for Crystalline-W which needs to transfer the parent's era in `slow_path()`, and the parent object can already end up being retired. We use the following trick. When the parent node is still *not* retired, we simply retrieve the birth era from the node. Otherwise, we retrieve REFS' value of the minimum era.

```

1 // Redefine RNODE to encode REFS links: steal one bit
2 // to indicate REFS nodes (also applies to other
3 // functions that previously used dummy RNODE)
4 #define IS_RNODE(x) (x & 0x1) // Check if a REFS link
5 #define RNODE(x) (x ^ 0x1) // Encode or decode REFS
6
7 // Another huge addend for the slow path
8 // (in addition to previously defined REFC_PROTECT)
9 const uint64 REFC_PROTECT_HANOVER = 1 << 62;
10
11 // Adds a special REFS-terminal node and list tainting
12 void traverse(Node* next) {
13     while (next != nullptr) {
14         Node* curr = next;
15         if (IS_RNODE(curr)) { // REFS-terminal node
16             // It is always the last node, exit
17             Node* refs = RNODE(curr);
18             if (FAA(&refs->refc, -1) == 1) free_batch(refs);
19             break;
20         }
21         next = SWAP(&curr->next, invptr); // Tainting
22         Node* refs = curr->blink;
23         if (FAA(&refs->refc, -1) == 1) free_batch(refs);
24     }
25
26 // Increments the global era, replaces regular FAA
27 // (needs to help other threads first)
28 void increment_era() {
29     if (slow_counter != 0) {
30         for (int i = 0; i < MAX_THREADS; i++) {
31             for (int j = 0; j < MAX_IDX; j++) {
32                 if (state[i][j].result.V == invptr)
33                     help_thread(i, j);
34             }
35         }
36         FAA(&global_era, 1);
37     }
38 }

```

```

37 void try_retire() { // This replacement is wait-free
38     uint64 min_birth = batch.refs->birth;
39     Node* last = batch.first;
40     // Also check odd tags to bound slow-path loops
41     for (int i = 0; i < MAX_THREADS; i++) {
42         for (int j = 0; j < MAX_IDX+2; j++) {
43             if (rsrv[i][j].list.V == invptr ||
44                 (rsrv[i][j].list.T & 0x1)) continue;
45             if (rsrv[i][j].era.V < min_birth ||
46                 (rsrv[i][j].era.T & 0x1)) continue;
47             if (last == batch.refs)
48                 return; // Ran out of nodes, exit
49             last->slot = &rsrv[i][j];
50             last = last->bnext;
51         }
52     }
53 // Retire, make it wait-free by list tainting
54 Node* curr = batch.first;
55 int64 cnt = -REFC_PROTECT;
56 for (; curr != last; curr = curr->bnext) {
57     Reservation* slot = curr->slot;
58     if (slot->list.V == invptr) continue;
59     Node* prev = SWAP(&slot->list.V, curr);
60     if (prev != nullptr) {
61         if (prev == invptr) { // Inactive previously
62             if (CAS(&slot->list.V, curr, invptr))
63                 continue; // Try to rollback
64         } else { // Tainted: traverse the chopped tail
65             if (!CAS(&curr->next, nullptr, prev))
66                 traverse(prev);
67         }
68     }
69     cnt++;
70 }
71 if (FAA(&batch.refs->refc, cnt) == -cnt)
72     free_batch(batch.refs);
73 batch.first = nullptr; batch.counter = 0;
74 }

```

Figure 2: Crystalline-W's try\_retire(), traverse(), and increment\_era().

```

1 // Hand over the parent object if it is retired
2 void handover_parent(Node* parent) {
3     if (parent && parent->blink != nullptr) {
4         Node* refs = get_refs_node(parent);
5         FAA(&refs->refc, REFC_PROTECT_HANOVER);
6         int64 cnt = -REFC_PROTECT_HANOVER;
7         for (int i = 0; i < MAX_THREADS; i++)
8             if (CAS(&parents[i], parent, nullptr)) cnt++;
9         FAA(&refs->refc, cnt);
10    }
11
12 uint64 get_birth_era(Node* node) { // Get parent's
13     if (node == nullptr) return 0; // birth era
14     uint64 birth_era = node->birth;
15     Node* link = node->blink;
16     // For already retired SLOT nodes, use REFS' value
17     if (link != nullptr && !IS_RNODE(link))
18         birth_era = link->birth;
19     return birth_era;
20 }

```

```

21 // Makes the tag+1 transition and detaches an old list
22 void detach_nodes(int i, int j, int tag) {
23     // A simple era tag transition: tag -> tag+1
24     CAS(&rsrv[i][j].era.T, tag, tag+1);
25     // Detach nodes and increment the list tag
26     do { // Bounded by MAX_THREADS (try_retire checks
27         old = rsrv[i][j].list; // the era tag
28         if (old.T != tag) break;
29         bool success = WCAS(&rsrv[i][j].list,
30                             old, { nullptr, tag+1 });
31     } while (!success);
32     return success ? old.V : invptr; // Previous value
33 }
34
35 // Get REFS node from any node in a batch
36 Node* get_refs_node(Node* node) {
37     Node* refs = node->blink;
38     if (IS_RNODE(refs)) refs = node; // This node itself
39     return refs;
40 }

```

Figure 3: Crystalline-W's utility functions for the slow path.

```

1 void slow_path(Node* obj, int index, Node* parent) {
2     // Getting parent's birth is tricky: for non-retir-
3     // ed nodes use 'birth', else retrieve the minimum
4     // birth from REFS, see get_birth_era() for details
5     uint64 parent_birth = get_birth_era(parent);
6     FAA(&slow_counter, 1);
7     state[TID][index].obj = obj;
8     state[TID][index].parent = parent;
9     state[TID][index].era = parent_birth;
10    uint64 tag = rsrv[TID][index].era.T;
11    state[TID][index].result = { invptr, tag };
12    uint64 prev_era = rsrv[TID][index].era.V;
13    do { // Bounded by MAX_THREADS
14        Node* list, * ptr = *obj;
15        uint64 curr_era = global_era;
16        if (curr_era == prev_era &&
17            WCAS(&state[TID][index].result,
18                { invptr, tag }, { nullptr, 0 })) {
19            rsrv[TID][index].era.T = tag+2;
20            rsrv[TID][index].list.T = tag+2;
21            FAA(&slow_counter, -1);
22            return ptr; // DONE
23        }
24        // Dereference previous nodes and update the era
25        if (rsrv[TID][index].list.V != nullptr) {
26            list = SWAP(&rsrv[TID][index].list.V,
27                       nullptr);
28            if (rsrv[TID][index].list.T != tag)
29                goto produced; // Result was just produced
30            if (list != invptr) traverse(list);
31            curr_era = global_era;
32        }
33        // WCAS fails only when the result is produced
34        WCAS(&rsrv[TID][index].era,
35            { prev_era, tag }, { curr_era, tag });
36        prev_era = curr_era;
37    } while (state[TID][index].result.V == invptr);
38    list = detach_nodes(TID, index, tag); //tag+1 state
39    produced:
40    ptr = state[TID][index].result.V;
41    uint64 era = state[TID][index].result.T;
42    rsrv[TID][index].era.V = era;
43    rsrv[TID][index].era.T = tag+2;
44    rsrv[TID][index].list.T = tag+2;
45    // Check if the obtained node is already retired
46    if (ptr && ptr->blink != nullptr) {
47        Node* refs = get_refs_node(ptr);
48        FAA(&refs->refc, 1);
49        if (list != invptr) traverse(list);
50        list = SWAP(&rsrv[TID][index].list.V,
51                   RNODE(refs)); // Put a REFS-terminal node
52    }
53    FAA(&slow_counter, -1);
54    // Traverse the previously detached list
55    if (list != invptr) traverse(list);
56    // Hand over the parent to all helper threads
57    handover_parent(parent);
58    return ptr; // DONE
59 }

60 void help_thread(int i, int j) {
61     Tag<void> result = state[i][j].result;
62     if (result.V != invptr) return;
63     uint64 era = state[i][j].era;
64     Node* parent = state[i][j].parent;
65     if (parent != nullptr) {
66         rsrv[TID][MAX_IDX].list.V = nullptr;
67         rsrv[TID][MAX_IDX].era.V = era;
68         parents[TID] = parent; // Advertise for a handover
69     }
70     Node* obj = state[i][j].obj;
71     uint64 tag = rsrv[i][j].era.T;
72     if (tag != result.T) goto changed;
73     uint64 curr_era = global_era;
74     do { // Bounded by MAX_THREADS
75         prev_era = update_era(curr_era, MAX_IDX+1);
76         Node* ptr = *obj;
77         uint64 curr_era = global_era;
78         if (prev_era == curr_era) {
79             if (WCAS(&state[i][j].result, // Published the
80                     result, { ptr, curr_era })) { // result
81                 Node* list = detach_nodes(i, j, tag);
82                 if (list != invptr) traverse(list);
83                 do { // Set the new era, <= 2 iterations
84                     old = rsrv[TID][index].era;
85                     if (old.T != tag+1) break;
86                 } while (!WCAS(&rsrv[TID][index].era,
87                               old, { curr_era, tag+2}));
88                 // If the obtained node is already retired
89                 if (ptr && ptr->blink != nullptr) {
90                     Node* refs = get_refs_node(ptr);
91                     FAA(&refs->refc, 1);
92                     do { // Bounded by MAX_THREADS
93                         old = rsrv[TID][index].list;
94                         if (old.T != tag+1) break;
95                         ok = WCAS(&rsrv[TID][index].list,
96                                 old, { RNODE(refs), tag+2 }));
97                         if (ok && old.V != invptr) traverse(old.V);
98                         if (ok) goto done;
99                     } while (!ok);
100                    FAA(&refs->refc, -1); // Already inserted
101                } else { // A simple tag transition
102                    CAS(&rsrv[TID][index].list.V, tag+1, tag+2);
103                }
104                break;
105            }
106        } while (state[i][j].result == result);
107    done:
108    Node* lst = SWAP(&rsrv[TID][MAX_IDX+1].list.V, invptr)
109    traverse(lst);
110    changed: // If handover occurs, dereference the parent
111    if (parent != nullptr) {
112        if (SWAP(&parents[TID], nullptr) != parent) {
113            Node* refs = get_refs_node(parent);
114            if (FAA(&refs->refc, -1) == 1) free_batch(refs);
115        }
116        Node* lst = SWAP(&rsrv[TID][MAX_IDX].list.V, invptr)
117        traverse(lst);
118    } }

```

Figure 4: Crystalline-W's slow-path methods.

## References

- [1] R. Nikolaev and B. Ravindran. Universal Wait-Free Memory Reclamation. In *Proceedings of the 25th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, PPOPP '20, pages 130–143, New York, NY, USA, 2020. ACM.
- [2] R. Nikolaev and B. Ravindran. Snapshot-Free, Transparent, and Robust Memory Reclamation for Lock-Free Data Structures. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*, PLDI '21, pages 987–1002, New York, NY, USA, 2021. ACM.
- [3] P. Ramalhete and A. Correia. Hazard Eras - Non-Blocking Memory Reclamation (Full Version). 2017.