

# Wait Free Atomic Snapshot Implementations

Sharan Narasimhan, CS20MTECH14003  
November 11th, 2020

## 1. Foreword to Evaluators:

I have kept the code explanation concise for the sake of brevity. Please read the code comments for more specific explanations.

The experiments for mr\*w are written to experiments\_mr\*w.txt and the timeline of execution is written to snapshots\_file\_mr\*w.txt.

The helpers.cpp contains common methods used for both snapshots.

## 2. MRSW Snapshot

### Step 1:

The main() consists of the usual steps such as initialising files, getting input parameters, initialising the snapshot object, creating threads and running the process inside #pragma omp critical. This all is included in Line 165-210 inside main().

### Step 2:

N+1 threads are created, where the thread with id = n is the snapshot thread, which will execute the while loop to attain a clean snap k times, as seen in figure 1.

```

210 if (id == n) //snapshot collecting thread executes here
211 {
212     int clean_snap[n];
213
214     while (no_of_snapshots < k)
215     {
216         start_time = preprocess_timestamp(omp_get_wtime());
217
218         //scan
219         ss.scan(clean_snap);
220         time_now = preprocess_timestamp(omp_get_wtime());
221         write_to_file(snapshots_file, n, -1, clean_snap, time_now, -1, no_of_snapshots, true);
222         no_of_snapshots++;
223
224         //update avg and worst times
225         tot_time = preprocess_timestamp(omp_get_wtime()) - start_time;
226         avg_time += tot_time;
227
228         if (tot_time > worst_time)
229             worst_time = tot_time;
230
231         //wait for rand time
232         rand_time = snapshot_delay(generator);
233         usleep(rand_time);
234     }
235
236     stop_writing = true; //make other threads stop writing
237     experiments_file << "avg time: " << avg_time / k << ", worst time: " << worst_time << endl;
238 }

```

Figure 1

The snapshot thread obtains a clean snap, writes it to file, computes relevant statistics, and sleep for a random time.

### Step 3:

All the other writer threads will execute update() indefinitely till the snapshot thread sets stop\_writing to true. After this all threads exit the parallel section and the file streams are closed. All this is encapsulated in figure 2.

```

240     else //writing thread executes her
241     {
242         while (!stop_writing)
243         {
244
245             //update
246             rand_val = rand() % 100;
247             ss.update(id, rand_val);
248             time_now = preprocess_timestamp(omp_get_wtime());
249             write_to_file(snapshots_file, n, id, nullptr, time_now, rand_val, -1, false);
250
251             //wait for rand time
252             rand_time = writer_delay(generator);
253             usleep(rand_time);
254         }
255     }
256 }
257
258 ss.display();
259
260 // close files
261 input_file.close();
262 experiments_file.close();
263 snapshots_file.close();
264 }

```

Figure 2

### The snap\_value class:

Each atomic register contains an object of type snap\_value which looks like:

```

16  class snap_value
17  {
18
19  public:
20      int value;
21      int label;
22      int snap[MAX_THREADS];

```

Figure 3

### The mrsw\_snapshot\_obj class:

This class encapsulates the array of atomic<snap\_value> MRSW register and all the necessary functions such as collect(), update() and scan().

```

47  class mrsw_snapshot_obj
48  {
49  public:
50      int n_threads;
51      atomic<snap_value> s_table_snap_values[MAX_THREADS];
52
53      mrsw_snapshot_obj(int n)
54      {
55          n_threads = n;
56          snap_value dummy_sv(n_threads);
57
58          for (int i = 0; i < n_threads; i++)
59          {
60              s_table_snap_values[i].store(dummy_sv);
61          }
62      }

```

Figure 4

The update() function gets a clean snapshot, takes the old snap\_value object and creates the new entry to be written to s\_table\_snap\_values as seen below:

```

64  void update(int me, int new_value) // this is invoked after scan()
65  {
66
67      // get clean snap
68      int clean_snap[n_threads];
69      scan(clean_snap);
70
71      snap_value old_snap_value = s_table_snap_values[me].load();
72      snap_value new_sv(n_threads, clean_snap, new_value, old_snap_value.label + 1);
73
74      s_table_snap_values[me].store(new_sv);
75  }

```

Figure 5

The collect() function is fairly straightforward.

The scan() function follows the strategy, “return snap if clean collect is obtained, else steal snapshot of first thread that moved twice”. This logic can be seen below:

```

89     while (true)
90     {
91         while (true)
92         {
93             //get new collect
94             collect(new_copy);
95
96             for (int i = 0; i < n_threads; i++)
97             {
98                 if (old_copy[i] != new_copy[i])
99                 {
100                     if (moved[i])
101                     { //second move, get snap of this register
102                         cout << "found second move !" << endl;
103                         copy_b_to_a(clean_scan, s_table_snap_values[i].load().snap, n_threads);
104
105                         return;
106                     }
107                     else
108                     {
109                         moved[i] = true;
110                         // cout << "someone moved!" << endl;
111                         copy_b_to_a(old_copy, new_copy, n_threads); // old_copy = new_copy
112                         break;
113                     }
114                 }
115             }
116             copy_b_to_a(clean_scan, new_copy, n_threads);
117             return;
118         }
119     }
120 }

```

Figure 6

### 3. MRMW Snapshot

This algorithm is fairly similar to the previous MRSW algorithm except 1. Update() follows the strong freshness policy, 2. We use 2 objects in the MRMW class as explained below.

a) s\_table[m] contains objects of type mrmw\_entry as seen below:

```

19 class mrmw_entry
20 {
21
22 public:
23     int value;
24     int pid;
25     int seq_no;
26
27 mrmw_entry(int n) //construction for dummy snap value
28 {
29     value = -1;
30     pid = -1;
31     seq_no = -1;
32 }

```

Figure 7

b) help\_snaps[n] is an array of clean snaps corresponding to each thread. Its declaration can be seen in line 48 of the below figure.

```
42  class mrmw_snapshot_obj
43  {
44  public:
45      int n_threads;
46      int m_slots;
47      atomic<mrmw_entry> s_table[MAX_MRMW_ARRAY_SIZE];
48      atomic<array<int, MAX_MRMW_ARRAY_SIZE>> help_snaps[MAX_THREADS]; //array of n_thread atomic std::arrays
```

Figure 8

I was forced to use std::array instead of C++ arrays for this algorithm as the atomic<> type does not accept arrays.

The update() follows strong freshness and also includes the location of the register to update along with the value. Apart from that it is fairly straightforward.

The snap() also works the same as before whenever a double move is detected, the scanning thread will take the help\_snap of the array that moved twice. All this logic can be seen below:

```
106      while (true)
107      {
108          trythis:
109              while (true)
110              {
111                  collect_stdarray(new_copy);
112
113                  for (int i = 0; i < m_slots; i++)
114                  {
115                      if (old_copy[i] != new_copy[i])
116                      {
117                          //slot i in s_table[m_slots] has changed, get the pid of thread that moved
118                          thread_that_moved = s_table[i].load().pid;
119
120                          if (moved[thread_that_moved]) //second move, get snap of this register
121                          {
122                              auto help_snap_chosen = help_snaps[thread_that_moved].load();
123                              clean_scan = help_snap_chosen;
124                              return;
125                          }
126                          else
127                          {
128                              moved[thread_that_moved] = true;
129                              old_copy = new_copy;
130                              goto trythis;
131                          }
132                      }
133                  }
134                  clean_scan = new_copy;
135                  return;
136              }
137      }
138  }
```

Figure 9

#### 4. Correctness:

Please compile and run `mr*w_snapshot.cpp` and refer to the `snapshots_file_mr*w.txt` to confirm the correctness of the algorithm.

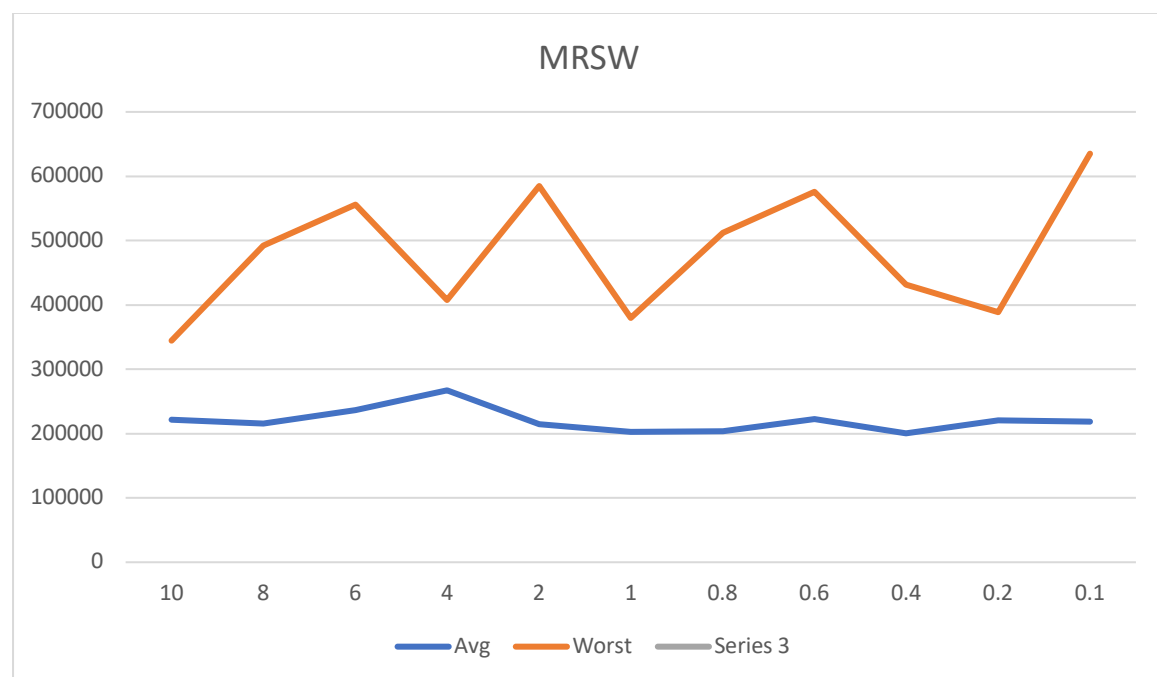
Please note that writing to these files was done sequentially by each thread to ensure that the timeline was in increasing order of time. Writes to files were done using the function `write_to_file()` in `helpers.cpp`. This will not affect the snapshot objects behaviour in any way, since writes are only done after an `update()` or a `scan()`.

#### 5. Analysis:

Please refer to `experiments_mr*w.txt` for all the data used for analysis.

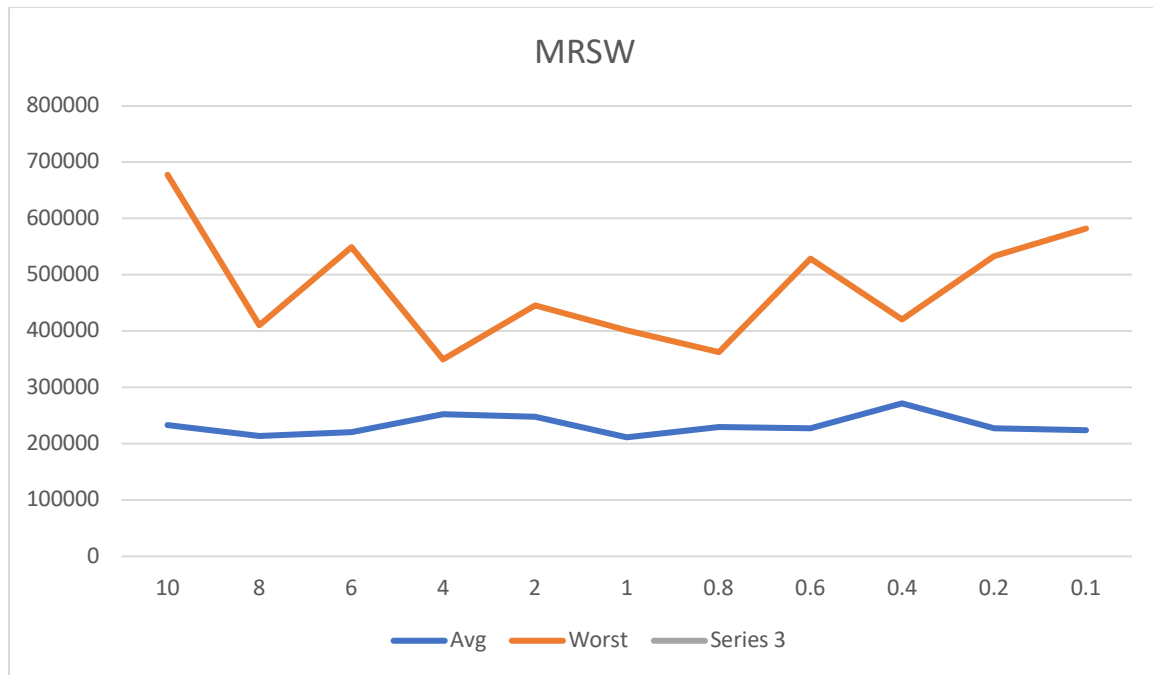
##### a) MRSW snapshot analysis:

The parameters were  $n = 10$ ,  $k = 50$  for all experiments. Each datapoint was obtained by averaging 5 iterations. Since it was not specified how to obtain the  $u_s/u_r$  ratio, I am keeping  $u_s$  set to 100 micro seconds and varying  $u_r$ .



##### a) MRMW snapshot analysis:

The parameters were  $n = 10$ ,  $m = 20$ ,  $k = 50$  for all experiments. Each datapoint was obtained by averaging 5 iterations. Since it was not specified how to obtain the  $u_s/u_r$  ratio, I am keeping  $u_s$  set to 100 micro seconds and varying  $u_r$ .



## Observations:

### 1. Both MRMW and MRSW have similar average times

This at first seems odd as we expected MRMW to experience more collisions since now collision occurs through multiple threads. However the probability that a register is updated is still  $1/n\_threads$  even in the case of MRMW. This is the same case for MRSW since at any given time the probability that the snapshot thread experiences a collision is  $1/n\_threads$  also.

Since #collisions is the primary factor affecting time taken for scan(), both the algos will perform similarly.

### 2. There is no relation between $u\_s/u\_w$ and avg time

This is a little odd and hard to explain. I have no concrete reasoning as to why this is occurring.

### 3. Worst time is randomly fluctuating:

This is expected, and we can see the "avg worst time" does not vary a lot.



