

# Wait-free Queues with Polylogarithmic Step Complexity

Hossein Naderibeni

supervised by Eric Ruppert

October 6, 2022

## Abstract

In this work, we are going to introduce a novel lock-free queue implementation. Linearizability and lock-freedom are standard requirements for designing shared data structures. All existing linearizable, lock-free queues in the literature have a common problem in their worst case called CAS Retry Problem. Our contribution is solving this problem while outperforming the previous algorithms.

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Related Work</b>	<b>5</b>
2.1	List-based Queues . . . . .	5
2.2	Universal Constructions . . . . .	7
2.3	Attia Fourier Lower Bound . . . . .	7
<b>3</b>	<b>Our Queue</b>	<b>8</b>
3.1	Pseudocode description . . . . .	15
3.2	Pseudocode . . . . .	18
<b>4</b>	<b>Proof of Correctness</b>	<b>23</b>
4.1	Basic Properties . . . . .	23
4.2	Ordering Operations . . . . .	27
4.3	Propagating Operations to the Root . . . . .	30

4.4	Correctness of GetEnqueue . . . . .	33
4.5	Correctness of IndexDequeue . . . . .	36
4.6	Linearizability . . . . .	39
<b>5</b>	<b>Analysis</b>	<b>42</b>
5.1	Garbage Collection or Getting rid of the infinite Arrays . . . . .	44
<b>6</b>	<b>Using Queues to Implement Vectors</b>	<b>45</b>
<b>7</b>	<b>Conclusion</b>	<b>46</b>

# 1 Introduction

Shared data structures have become an essential field in distributed algorithms research. We are reaching the physical limits of how many transistors we can place on a CPU core. The industry solution to provide more computational power is to increase the number of cores of the CPU. This is why distributed algorithms have become important. It is not hard to see why multiple processes cannot update sequential data structures designed for one process. For example, consider two processes trying to insert some values into a sequential linked list simultaneously. Processes  $p, q$  read the same tail node,  $p$  changes the next pointer of the tail node to its new node and after that  $q$  does the same. In this run,  $p$ 's update is overwritten. One solution is to use locks; whenever a process wants to do an update or query on a data structure, the process locks it, and others cannot use it until the lock is released. Using locks has some disadvantages; for example, one process might be slow, and holding a lock for a long time prevents other processes from progressing. Moreover, locks do not allow complete parallelism since only the one process holding the lock can make progress.

The question that may arise is, “What properties matter for a lock-free data structure?”, since executions on a shared data structure are different from sequential ones, the correctness conditions also differ. To prove a concurrent object works perfectly, we have to show it satisfies safety and progress conditions. A *safety condition* tells us that the data structure does not return wrong responses, and a *progress property* requires that operations eventually terminate.

The standard safety condition is called *linearizability*, which ensures that for any concurrent execution on a linearizable object, each operation should appear to take effect instantaneously at some moment between its invocation and response. Figure 1 is an example of an execution on a linearizable queue that is initially empty. The arrow shows time, and each rectangle shows the time between the invocation and the termination of an operation. Since `Enqueue(A)` and `Enqueue(B)` are concurrent, `Enqueue(B)` may or may not take effect before `Enqueue(A)`. The execution in Figure 2 is not linearizable since A has been enqueued before B, so it has to be dequeued first.

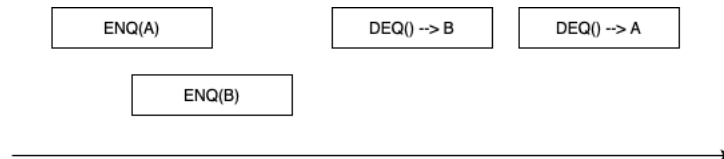


Figure 1: An example of a linearizable execution. Either `Enqueue(A)` or `Enqueue(B)` could take effect first since they are concurrent.

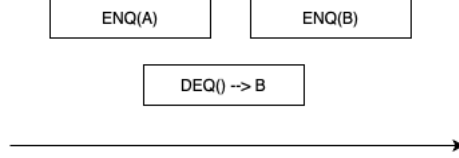


Figure 2: An example of an execution that is not linearizable. Since `Enqueue(A)` has completed before `Enqueue(B)` is invoked the `Deque()` should return A or nothing.

There are various progress properties; the strongest is wait-freedom, and the more common is lock-freedom. An algorithm is *wait-free* if each operation terminates after a finite number of its own steps. We call an algorithm *lock-free* if, after a sufficient number of steps, one operation terminates. A wait-free algorithm is also lock-free but not vice versa; in an infinite run of a lock-free algorithm there might be an operation that takes infinitely many steps but never terminates.

In section 2 we talk about previous queues and their common problems. We also talk about polylogarithmic construction of shared objects.

Jayanti [9] proved an  $\Omega(\log p)$  lower bound on the worst-case shared-access time complexity of  $p$ -process universal constructions. He also introduced [1] a construction that achieves  $O(\log^2 p)$  shared accesses. Here, we first introduce a universal construction using  $O(\log p)$  CAS operations [10]. In section 3 we introduce a polylogarithmic step wait-free universal construction. Our main ideas in of the universal construction also appear in our Queue Algorithm (??). The main short come of our universal construction is using big CAS objects. We use the universal construction as a stepping stone towards our queue algorithm, so we will not explain it in too much detail.

In section 4 we introduce a concurrent wait-free datastructure, to agree on the order of the operations invoked on some processes.

In section 5 we introduce our main work, the queue; prove its linearizability and wait-freeness.

## 2 Related Work

### 2.1 List-based Queues

In the following paragraphs, we look at previous lock-free queues. Michael and Scott [13] introduced a lock-free queue which we refer to as the MS-queue. A version of it is included in the standard Java Concurrency Package. Their idea is to store the queue elements in a singly-linked list (see Figure 3). Head points to the first node in the linked list that has not been dequeued, and Tail points to the last element in the queue. To insert a node into the linked list, they use atomic primitive operations like LL/SC or CAS. If  $p$  processes try to enqueue simultaneously, only one can succeed, and the others have to retry. This makes the amortized number of steps to be  $\Omega(p)$  per enqueue. Similarly, dequeue can take  $\Omega(p)$  steps.

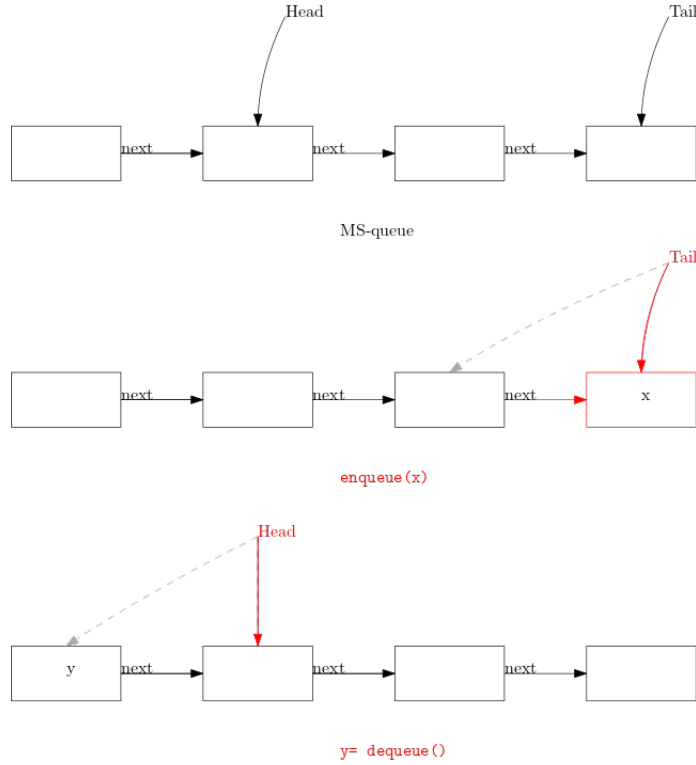


Figure 3: MS-queue structure, enqueue and dequeue operations. In the first diagram the first element has been dequeued. Red arrows show new pointers and gray dashed arrows show the old pointers.

Moir, Nussbaum, and Shalev [14] presented a more sophisticated queue by using the elimination technique. The elimination mechanism has the dual purpose of allowing operations to complete in parallel and reducing contention for the queue. An Elimination Queue consists of an MS-queue augmented with an elimination array. Elimination works by allowing opposing pairs of concurrent operations such as an enqueue and a

dequeue to exchange values when the queue is empty or when concurrent operations can be linearized to empty the queue. Their algorithm makes it possible for long-running operations to eliminate an opposing operation. The empirical evaluation showed the throughput of their work is better than the MS-queue, but the worst case is still the same; in case there are  $p$  concurrent enqueues, their algorithm is not better than MS-queue.

Hoffman, Shalev, and Shavit [8] tried to make the MS-queue more parallel by introducing the Baskets Queue. Their idea is to allow more parallelism by treating the simultaneous enqueue operations as a basket. Each basket has a time interval in which all its nodes' enqueue operations overlap. Since the operations in a basket are concurrent, we can order them in any way. Enqueues in a basket try to find their order in the basket one by one by using **CAS** operations. However, like the previous algorithms, if there are still  $p$  concurrent enqueue operations in a basket, the amortized step complexity remains  $\Omega(p)$  per operation.

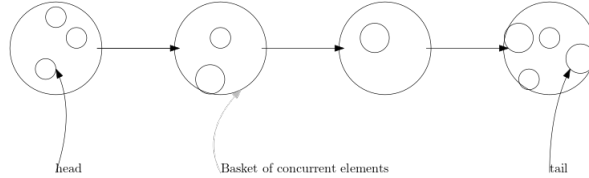


Figure 4: Baskets queue idea. There is a time that all operations in a basket were running concurrently, but only one has succeeded to do **CAS**. To order the operations in a basket, the mechanism in the algorithm for processes is to **CAS** again. The successful process will be the next one in the basket and so on.

Ladan-Mozes and Shavit [12] presented an Optimistic Approach to Lock-Free FIFO Queues. They use a doubly-linked list and do fewer **CAS** operations than MS-queue. But as before, the worst case is when there are  $p$  concurrent enqueues which have to be enqueued one by one. The amortized worst-case complexity is still  $\Omega(p)$  **CASes**.

Hendler et al. [6] proposed a new paradigm called flat combining. Their queue is linearizable but not lock-free. Their main idea is that with knowledge of all the history of operations, it might be possible to answer queries faster than doing them one by one. In our work we also maintain the whole history. They present experiments that show their algorithm performs well in some situations.

Gidenstam, Sundell, and Tsigas [4] introduced a new algorithm using a linked list of arrays. Global head and tail pointers point to arrays containing the first and last elements in the queue. Global pointers are up to date, but head and tail pointers may be behind in time. An enqueue or a dequeue searches in the head array or tail array to find the first unmarked element or last written element (see Figure 5). Their data

structure is lock-free. Still, if the head array is empty and  $p$  processes try to enqueue simultaneously, the step complexity remains  $\Omega(p)$ .

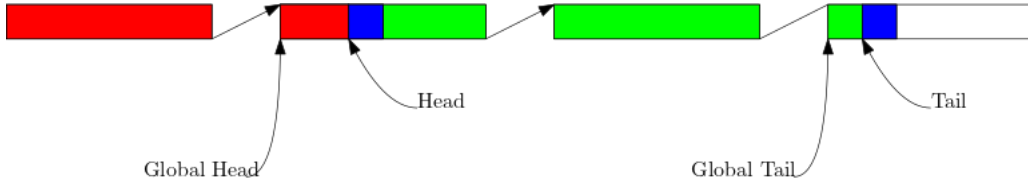


Figure 5: Global pointers point to arrays. Head and Tail elements are blue, dequeued elements are red and current elements of the queue are green.

Kogan and Petrank [11] introduced wait-free queues based on the MS-queue and use Herlihy’s helping technique to achieve wait-freedom. Their step complexity is  $\Omega(p)$  because of the helping mechanism.

In the worst-case step complexity of all the list-based queues discussed above, there is a  $p$  term that comes from the case all  $p$  processes try to do an enqueue simultaneously. Morrison and Afek call this the *CAS retry problem* [15]. It is not limited to list-based queues and array-based queues share the CAS retry problem as well [17, 16, 2]. We are focusing on seeing if we can implement a queue in sublinear steps in terms of  $p$  or not.

## 2.2 Universal Constructions

Herlihy discussed the possibility of implementing shared objects from other objects [7]. A *universal construction* is an algorithm that can implement a shared version of any given sequential object. We can implement a concurrent queue using a universal construction. Jayanti proved an  $\Omega(\log p)$  lower bound on the worst-case shared-access time complexity of  $p$ -process universal constructions [9]. He also introduced a construction that achieves  $O(\log^2 p)$  shared accesses [1]. His universal construction can be used to create any data structure, but its implementation is not practical because of using unreasonably large-sized CAS operations.

Ellen and Woelfel introduced an implementation of a Fetch&Inc object with step complexity of  $O(\log p)$  using  $O(\log n)$ -bit LL/SC objects, where  $n$  is the number of operations [3]. Their idea has similarities to Jayanti’s construction, and they represent the value of the Fetch&Inc using the history of successful operations.

## 2.3 Attiya Fourier Lower Bound

### 3 Our Queue

Jayanti and Petrovic introduced a wait-free polylogarithmic multi-enqueuer single-dequeueer queue [10]. We benefit from some ideas of their work to design a polylogarithmic multi-enqueuer multi-dequeueer queue. Our algorithm despite them does not use **CAS** operations with big words and does not put a limit on the number of concurrent operations. In our model there are  $p$  processes doing **Enqueue()**, **Dequeue()** operations concurrently. We use a shared tree among the processes (see Figure 6) to agree on one total ordering on the operations invoked by processes. Each process has a leaf which the order of operations invoked by the process is stored in it. When a process wishes to do an operation it appends the operation to its leaf and then tries to propagate its new operation up to the tree's root. In each node the ordering of operations propagated up to it is stored. All processes agree on the sequence stored in the root and it is defined to be the linearization ordering.

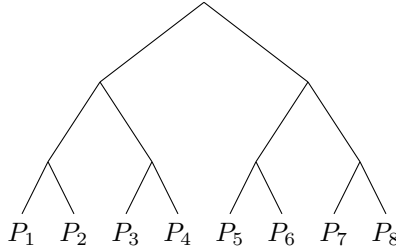


Figure 6: Each process has a leaf and in each node there is an ordering of operations stored. Each node tries to propagate its operations up to the root, which stores the total ordering of all operations.

#### *Add sequence to nodes*

We could implement the sequence stored in each node using an array of the queue operations and append some operations to the sequence by doing **k-CAS** operation on the end of the array. To do a propagate step on node  $n$  in the tree, we aggregate the operations from node  $n$ 's both children (that have not already been propagated to  $n$ ) and try to append them into  $n$ . We call this procedure **REFRESH( $n$ )**. The main idea is that if we call **REFRESH( $n$ )** twice, the operations in  $n$ 's children before the first **REFRESH( $n$ )** are guaranteed to be in  $n$ . Because if both of the **REFRESH( $n$ )**s fail to do **k-CAS** then there is another instance of **Refresh** in between which has succeeded to do **CAS** and has already appended the operations that the first **Refresh** was trying to append. This mechanism makes us overcome the **CAS** Retry Problem.

*fix ||*



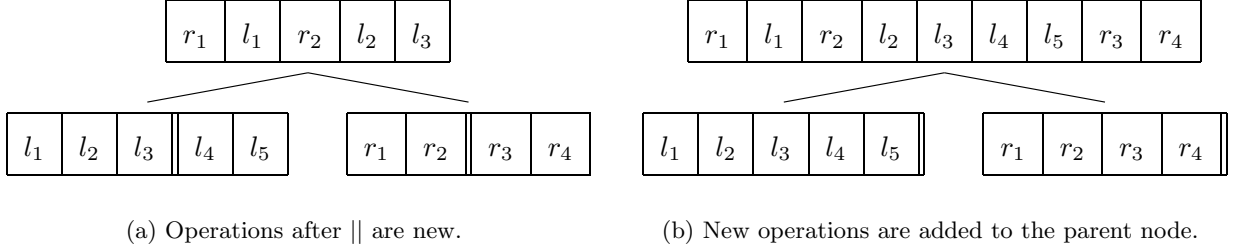


Figure 7: Before and after of a  $\text{REFRESH}(n)$  with successful CAS. Operations propagating from the left child are numbered with  $l$  and from the right child by  $r$  and the operations in children after  $||$  are new.

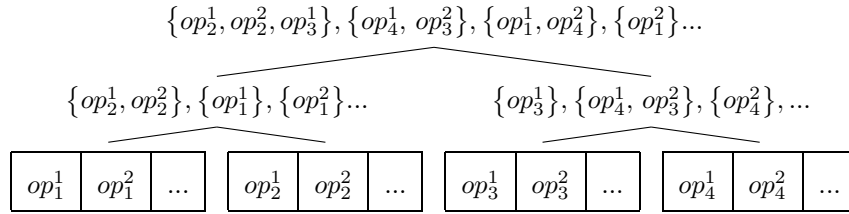


Figure 8: In each internal node, we store the set of all the operations propagated together, and one can arbitrarily linearize the sets of concurrent operations among themselves. Since we linearize operations when they are added to the root, ordering the blocks in the root is important.

The solution for implementing the orderings in the tree told above is not efficient, because there are big CASes and operations information are copied all the way up to the root. Instead of storing operations explicitly in the nodes, we can keep track of some statistics of them. This allows us to CAS fixed-size objects in each  $\text{REFRESH}(n)$ . To do that, we introduce blocks that only contain the number of operations from the left and the right child in a **Refresh()** procedure and only propagate the statistics block of the new operations. In each **Refresh** there is at most one operation from each process trying to be propagated, because one operation cannot invoke two operations concurrently. Furthermore since the operations in a  $\text{REFRESH}$  step are concurrent we can linearize them among themselves in any order we wish. Note that if two operations are in read one  $\text{REFRESH}$  step in a node they are going to be propagated up to the root together. Our choice is to put the operations propagated from the left child before the operations propagated from the right child. In this way if we know the number of operations from the left child and the number of operations from the right child in a block we have a complete ordering on the operations.

A process may wish to know the  $i$ th propagated operation or the rank of a propagated operation in the

linearization. In our case of implementing a queue, we can make an assumption that one process only wishes to know the rank of a dequeue and one tries to get an enqueue with its rank. **enqueues** and **dequeues** are appended to the tree and when we want to find the response to a **dequeue**, we compute the place of the dequeue in the linearization and using the rank of the dequeue among dequeues and some information stored in the root we compute which enqueue is the answer to the dequeue or if the answer is null. If the answer was some enqueue we find the enqueue using **DSearch(i)** and **GetEnqueue(n,b,i)**. **DSearch(i)** finds the block containing the  $i$ th enqueue in the root and **GetEnqueue(n,b,i)** finds its sub-block recursively to reach a leaf. **Index()** is similar but more complicated, finding super-blocks from a leaf to the root. The main challenge in each level of **Get(i)** and **Index(op)** is that it should take polylogarithmic steps with respect to  $p$ . After appending operation **op** to the root, processes can find out information about the linearization ordering using **Get(i)** and **Index(op)**. Each block stores an extra constant amount of information (like prefix sums) to allow binary searches to find the required block in a node quickly.

**Implementing Queue using Block Tree** In this work, we design a queue with  $O(\log^2 p + \log n)$  steps per operation, where  $n$  is the number of total operations invoked. We avoid the  $\Omega(p)$  worst-case step complexity of existing shared queues based on linked lists or arrays (CAS Retry Problem). A queue stores a sequence of elements and supports two operations, enqueue and dequeue. **Enqueue(e)** appends element **e** to the sequence stored. **Dequeue()** removes and returns the first element among in the sequence. If the queue is empty it returns **null**. Knowing index  $i$  is the tail of the queue, we can return the dequeue response using **Get(i)**. So in the rest we modify block tree to compute **i** for each **Dequeue()** to achieve a FIFO queue.

**GETINDEX(i)** returns the  $i$ th operation stored in the block tree sequence. We do that by finding the block  $b_i$  containing  $i$ th element in the root, and then recursively finding the subblock of  $b_i$  which contains  $i$ th element. To make this recursive search faster, instead of iterating over all elements in sequence of blocks we store prefix sum of number of elements in the blocks sequence and pointers to make **BinarySearch** faster.

Furthermore, in each block, we store the prefix sum of left and right elements. Moreover, for each block, we store two pointers to the last left and right subblock of it (see fig 11 and 10).

Starting from the root, **GETINDEX(i)** **BinarySearches**  $i$  in the prefix sum array to find block containing  $i$ th operation, then continues recursively calling **GETELEMENT(b,i)** to find  $i$ th element of block  $b$ . From lemma 28 we know a block size is at most  $p$ . So **BinarySearch** takes at most  $(O)(\log p)$ , since with knowing pointers

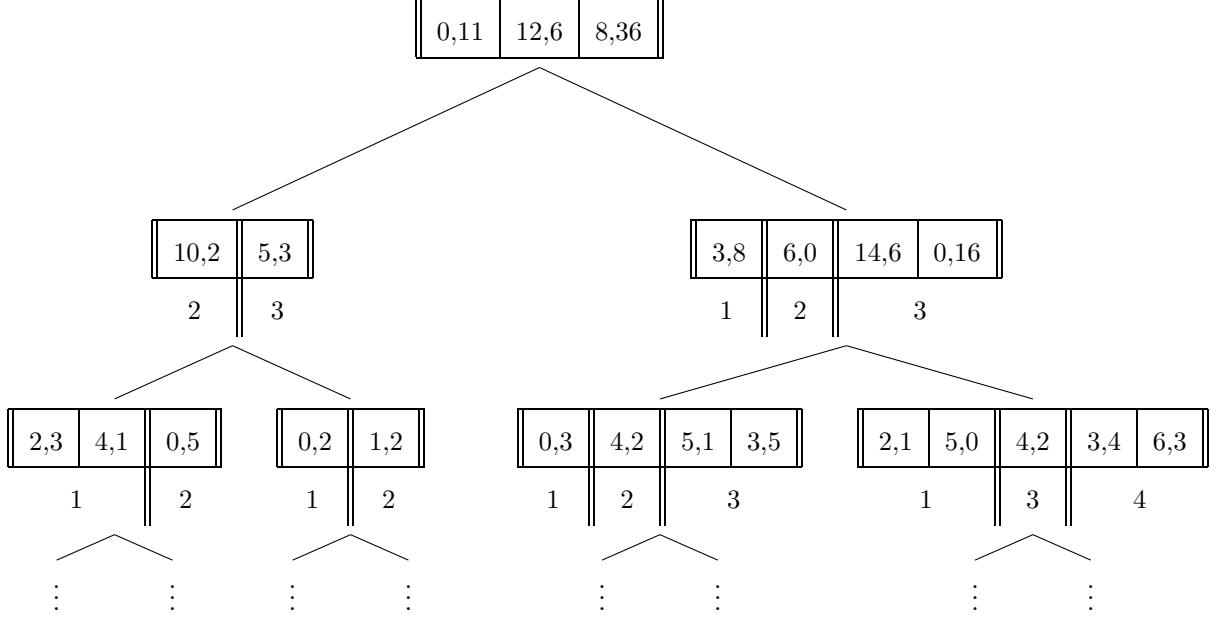


Figure 9: Showing concurrent operation sets with blocks. Each block consists of a pair(left, right) indicating the number of operations from the left and the right child, respectively. Block (12,6) in the root contains blocks (10,2) from the left child and (6,0) from the right child. Blocks between two lines || are propagated together to the parent. For example, Blocks (2,3) and (4,1) from the leftmost leaf and (0,2) from its sibling are propagated together into the block (10,2) in their parent. The number underneath a group of blocks in a node indicates which block in the node's parent those blocks were propagated to. Each block  $b$  in node  $n$  is the aggregation of blocks in the children of  $n$  that are newly read by the `PROPAGATE()` step that created block  $b$ . For example, the third block in the root (8,36) is created by merging block (5,3) from the left child and (14,6) and (0,16) from the right child. Block (5,3) also points to elements from blocks (0,5) and (1,2). We choose to linearize operations in a block from the left child before those from the right child as a convention. Operations within a block of the root can be ordered in any way that is convenient. In effect, this means that if there are concurrent new blocks in a `REFRESH()` step from several processes we linearize them in the order of their process ids. So for example operations aggregated in block (10,2) are in the order (2,3),(4,1),(0,2). All blocks from the left child will come before the right child and the order of blocks of each child is preserved among themselves.

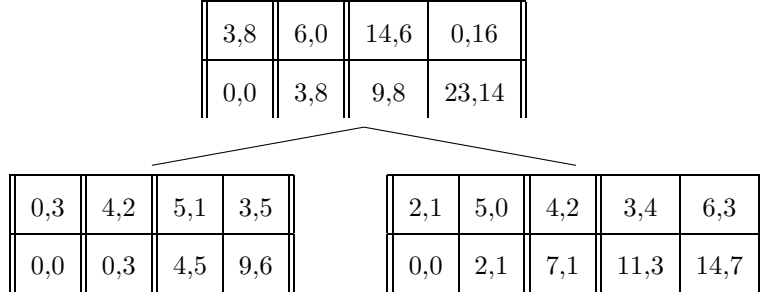


Figure 10: Using Prefix sums in blocks. When we want to find block  $b$  elements in its children, we can use binary search. The number below each block shows the count of elements in the previous blocks.

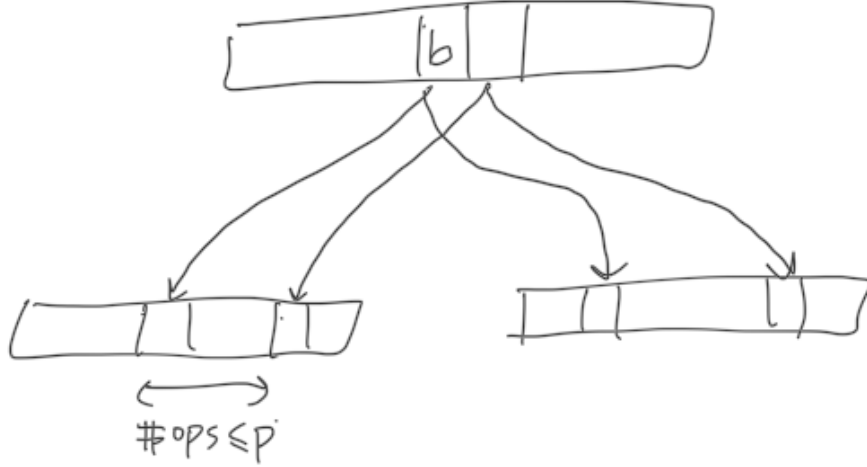


Figure 11: Block have pointers to the starting block of theirs for each child.

of a block and its previous block we can determine the base (domain ?) to search and its size is  $O(p)$ .

**CreateBlock** `CreateBlock(n)` returns a block containing new operations of  $n$ 's children.  $b'.end_{left}$  stores the index of the rightmost subblock of left child of  $b$ 's previous block. Other attributes are assigned values followed by definition.

**Computing** `Get(n, b, i)`

**How** `Refresh(n)` works.

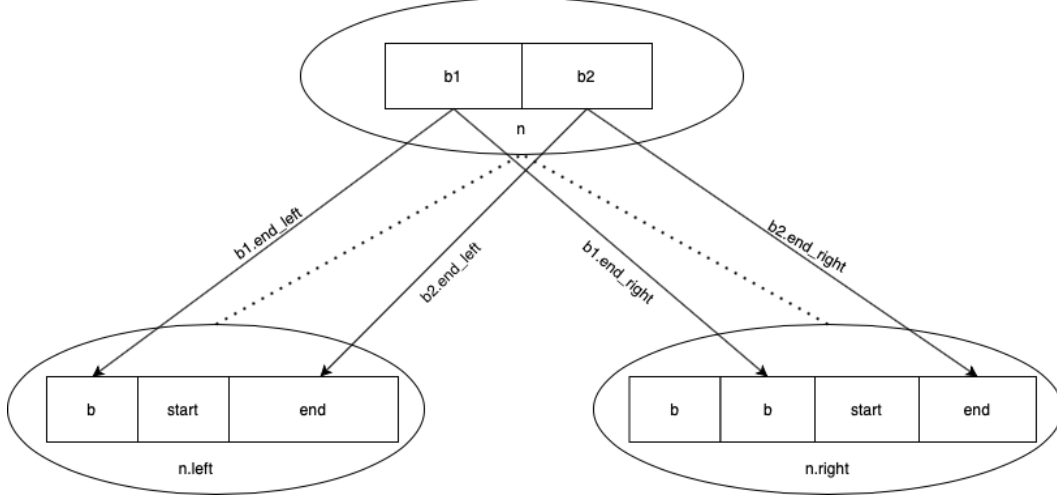


Figure 12: Snapshot of a CreateBlock

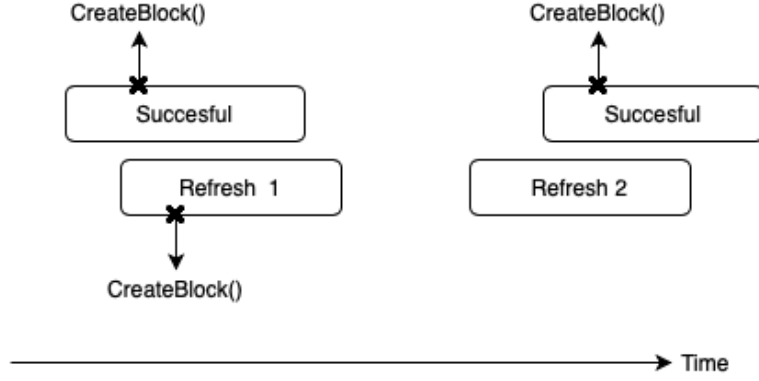


Figure 13: The second failed Refresh is assuredly concurrent to a Successful Refresh with CreateBlock line after first failed Refresh's CreateBlock.

**Computing superblock** Let  $i$  be the value  $R_n$ , a successful instance of **Refresh** on node  $n$  reads from  $n.head$ .  $R_n$  does a successful  $CAS(\text{null}, b)$  into  $n.blocks[i]$ . Let  $p$  be  $n.parent$ . Without loss of generality for the rest of this section assume  $n$  is the left child of  $p$ . From Lemma 27 we know there could be only one  $p.Refresh$  propagating  $b$ . Let  $R_p$  be the first successful  $p.Refresh$  that reads some value greater than  $i$  for  $left.head$  and contains  $b$  in its created block in Line 317. Let the index of the block  $R_p$  put in  $p.blocks$  be  $j$ .

Since the index of the superblock of  $b$  is not known until  $b$  is propagated,  $R_n$  cannot set the **super** field of  $b$  while creating it. One approach is to set the **super** field of  $b$  by  $R_n$  after propagating  $b$  to  $p$ . This solution would not be efficient because there might be  $p$  subblocks in the block  $R_p$  propagated needing to update the

`super` field. However intuitively, once  $b$  is installed, its superblock is going to be close to  $n.\text{parent.head}$  at the time of installation. One idea is that if we know the approximate position of the superblock of  $b$  then we can search for the real superblock when we wished to know the superblock of  $b$  i.e.  $b.\text{super}$  does not have to be the exact location of the superblock of  $b$ , but we want it to be close to  $j$ . We can set  $b.\text{super}$  to  $n.\text{parent.head}$  while creating  $b$ , but the problem is that there might be many  $p.\text{Refreshes}$  that could happen after reading  $p.\text{head}$  by  $R_n$  and before propagating  $b$  to  $p$ . If we set  $b.\text{super}$  to  $p.\text{head}$  after appending  $b$  to  $n.\text{blocks}$  (Line 326),  $R_n$  might go to sleep at some time after installing  $b$  and before setting  $b.\text{super}$ . In this case the next  $\text{Refreshes}$  on  $n$  and  $n.\text{parent}$  help fill in the value of  $b.\text{super}$ .

Block  $b$  is appended to  $n.\text{blocks}[h]$  on Line 320. After appending  $b$ ,  $b.\text{super}$  is set on Line 326 of a call to `Advance` from  $n.\text{Refresh}$  by the same process or another process or maybe an  $n.\text{parent.Refresh}$ . We want to bound how far  $b.\text{super}$  is from the index of  $b$ 's superblock, which is created by a successful  $n.\text{parent.Refresh}$  that propagates  $b$ .

**Queue from tree** Now, we describe how to use the tree to implement a queue. Consider the following execution of operations. `Enqueue(e)` appends an operation with input argument  $e$  in the block tree. What should a `Dequeue()` return? To compute the response of a `Dequeue()`, process  $p$  first appends a `DEQ` operation to the tree. Then  $p$  finds the rank of the `DEQ` using `Index()`, the rank of the `DEQ` and the information stored in the root about the queue  $p$  computes the rank of the `ENQ` having the answer of the `DEQ`. Finally  $p$  returns the argument of that `ENQ` using `Get(i)`.

ENQ(5)	ENQ(2)	DEQ()	ENQ(3)	DEQ()	DEQ()	DEQ()	ENQ(4)	ENQ(6)	DEQ()
--------	--------	-------	--------	-------	-------	-------	--------	--------	-------

Table 1: An example histoy of operations on the queue

A non-null dequeue is one that returns a non-null value. In the example above, `Dequeue()` operations return 5, 2, 3, null, 4 in order. Before `ENQ(4)` the queue gets empty so the last `DEQ()` returns null. If the queue is non-empty and  $r$  `Dequeue()` operations have returned a non-null response, then  $i$ th `Dequeue()` returns the input of the  $r + 1$ th `Enqueue()`. So, in order to answer a `Dequeue`, it's sufficent to know the size of the queue and the number of previous non-null dequeues.

In the Block Tree, we did not store the sequence of operations explicitly but instead stored blocks of concurrent operations to optimize `Propagate` steps and increase parallelism. So now the problem is to find

the result of each Dequeue. From lemma 28 we know we can linearize operations in a block in any order; here, we choose to decide to put Enqueue operations in a block before Dequeue operations. In the next example, operations in a cell are concurrent. `DEQ()` operations return `null`, 5, 2, 1, 3, 4, `null` respectively. We will next describe how these values can be computed efficiently.

<code>DEQ()</code>	<code>ENQ(5), ENQ(2), ENQ(1), DEQ()</code>	<code>ENQ(3), DEQ()</code>	<code>ENQ(4), DEQ(), DEQ(), DEQ(), DEQ()</code>
--------------------	--	----------------------------	---

Table 2: An example history of operation blocks on the queue

Now, we claimed that by knowing the current size of the queue and the number of non-null dequeue operations before the current dequeue, we could compute the index of the resulting `Enqueue()`. We apply this approach to blocks; if we store the size of the queue after each block of operations happens and the number of non-null dequeues dequeues till a block, we can compute each dequeue's index of result in  $O(1)$  steps.

	<code>DEQ()</code>	<code>ENQ(5), ENQ(2), ENQ(1), DEQ()</code>	<code>ENQ(3), DEQ()</code>	<code>ENQ(4), DEQ(), DEQ(), DEQ(), DEQ()</code>
<code>#enqueues</code>	0	3	1	1
<code>#dequeues</code>	1	1	1	4
<code>#non-null dequeues</code>	0	1	2	5
<code>size</code>	0	2	2	0

Table 3: Augmented history of operation blocks on the queue

Size and the number of non-null dequeues for  $b$ th block could be computed this way:

`size[b] = max(size[b-1] +enqueues[b] -dequeues[b], 0)`

`non-null dequeues[b] = non-null dequeues[b-1] +dequeues[b] -size[b-1] -enqueues[b]`

Given `DEQ` is in block  $b$ , `response(DEQ)` would be:

`(size[b-1] - index of DEQ in the block's dequeus >= 0) ? ENQ[non-null dequeues[b-1] + index of DEQ in the block's dequeus] : null;`

### 3.1 Pseudocode description

**Specification** A Queue is a shared data structure that stores a sequence of elements. It has two methods `Enqueue(e)` and `Dequeue()`. `Enqueue(e)` adds `e` to the end of the sequence. `Dequeue()` returns the first

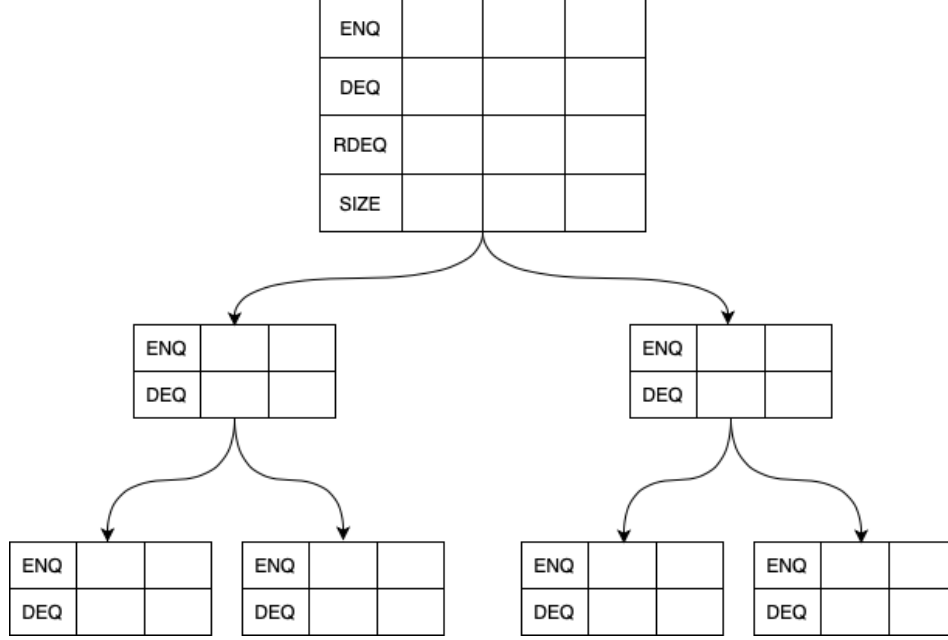


Figure 14: Fields stored in the Queue nodes.

element stored in the sequence and removes it from the sequence.

**Tree** In order to reach an agreement on the order of operations among  $p$  processes, we use a Tournament Tree. Leaf  $l_i$  is assigned to a process  $i$ . Each process adds  $op$  to its leaf. In each internal node an ordering of operations in its subtree is stored. All processes agree on the total ordering of all operations stored in the root. This ordering will be the linearization of the operations.

**Implicit Storing Blocks** For efficiency, instead of storing explicit sequence of operations in nodes of the Tournament Tree, we use Blocks. A Block is a constant size object that implicitly represents a sequence of operations. In each node there is an array of Blocks.

Block  $b$  contains subblocks in the left and right children. WLOG left subblocks of  $b$  are some consecutive blocks in the left child starting from where previous block of  $b$  has ended to the the end of  $b$ . See Figure 12 .

We store ordering among **operations** in the tournament tree constructed by **nodes**. In each **node** we store pointers to its relatives, an array of **blocks** and an index to the first empty **block**. Furthermore in **leaf** nodes there is an array of **operations** where each **operation** is stored in one cell with the same index in **blocks**. There is a **counter** in each **node** incrementing after a successful **Refresh** step. It means after that some bunch of **blocks** in a node have propagated into the parent then the **counter** increases. Each new



block added to a node sets its **time** regarding **counter**. This helps us to know which blocks have aggregated together to a block, not precisely though. We also store the index of the aggregated **block** of a **block** with **time**  $i$  in **super**[ $i$ ].

In each **block** we store 4 essential stats that implicitly summarize which operations are in the block **num<sub>enq-left</sub>**, **num<sub>deq-left</sub>**, **num<sub>enq-right</sub>**, **num<sub>deq-right</sub>**. In order to make **BinarySearch()**es faster we store prefix sums as well and there are some more general stats that help to make pseudocode more readable but not necessary.

To compute the head of the **queue** before a **dequeue** two more fields are stored in the root **size** and **sum<sub>non-null deq</sub>**. **size** in a **block** shows the number of elements after the **block** has finished and **sum<sub>non-null deq</sub>** is the total number of non-null dequeues till the **block**.

**Enqueue(e)** just appends an operation with element **e** to the root. **Dequeue()** appends an operation to the root and computes its ordering and the **enqueue** operation containing the head before it calling **ComputeHead()** and then **gets** and returns the operation's element.

**Append(op)** adds **op** to the invoking process's leaf's **ops** and **blocks**, propagates it up to the root and if the **op** is a dequeue returns its order in residing block in the root and the block's index. As we said later **Propagate** assuredly aggregates new blocks to a block in the parent by calling **Refresh** two times. **Refresh(n)** creates a block, tries to CAS it into the **pn**'s **blocks** and if it was successful updates **super** and **counter** in both of **n**'s children.

We only want to know the **element** of **enqueue** operations and compute ordering for **dequeue** operations. That's the reason here **Get()** searches between enqueues only and **Index()** returns ordering of a dequeue among dequeues. **Get(n, b, i)** decides the requested element is in which child of **n** and continues to search recursively. **index(n, i, b)** calculates the ordering of the given operation in **n**'s parent each step and finally returns the result among total ordering.

## 3.2 Pseudocode

---

### Algorithm Tree Fields Description

---

#### ◇ *Shared*

- A binary tree of Nodes with one leaf for each process. root is the root node.

#### ◇ *Local*

- *Node* leaf: process's leaf in the tree.

#### ► *Node*

- *\*Node* left, right, parent : Initialized when creating the tree.
- *Block[]* blocks : Initially blocks[0] contains an empty block with all fields equal to 0.
- *int* head= 1: #blocks in blocks. blocks[0] is a block with all integer fields equal to zero.

#### ► *Block*

- *int* super : approximate index of the superblock, read from parent.head when appending the block to the node

#### ► *RootBlock* extends *InternalBlock*

- *int* size : size of the queue after performing all operations in the prefix for this block

#### ► *InternalBlock* extends *Block*

- *int* end<sub>left</sub>, end<sub>right</sub> : indices of the last subblock of the block in the left and right child
- *int* sum<sub>enq-left</sub> : # enqueues in left.blocks[1..end<sub>left</sub>]
- *int* sum<sub>deq-left</sub> : # dequeues in left.blocks[1..end<sub>left</sub>]
- *int* sum<sub>enq-right</sub> : # enqueues in right.blocks[1..end<sub>right</sub>]
- *int* sum<sub>deq-right</sub> : # dequeues in right.blocks[1..end<sub>right</sub>]

#### ► *LeafBlock* extends *Block*

- *Object* element : Each block in a leaf represents a single operation. If the operation is enqueue(x) then element=x, otherwise element=null.
- *int* sum<sub>enq</sub>, sum<sub>deq</sub> : # enqueue, dequeue operations in the prefix for the block

---

*Abbreviations used in the code and the proof of correctness.*

- $\text{blocks}[b].\text{sum}_x = \text{blocks}[b].\text{sum}_{x\text{-left}} + \text{blocks}[b].\text{sum}_{x\text{-right}}$  (for  $b \geq 0$  and  $x \in \{\text{enq}, \text{deq}\}$ )
- $\text{blocks}[b].\text{num}_x = \text{blocks}[b].\text{sum}_x - \text{blocks}[b-1].\text{sum}_x$   
(for  $b > 0$  and  $x \in \{\text{enq}, \text{deq}, \text{enq-left}, \text{enq-right}, \text{deq-left}, \text{deq-right}\}$ )

---

## Algorithm Queue

---

```
201: void Enqueue(Object e)                                ▷ Creates a block with element e and adds it to the tree.
202:     block newBlock= new(LeafBlock)
203:     newBlock.element= e
204:     newBlock.sumenq= leaf.blocks[leaf.head].sumenq+1
205:     newBlock.sumdeq= leaf.blocks[leaf.head].sumdeq
206:     leaf.Append(newBlock)
207: end Enqueue

208: Object Dequeue()                                       ▷ Creates a block with null value element, appends it to the tree and returns its response.
209:     block newBlock= new(LeafBlock)
210:     newBlock.element= null
211:     newBlock.sumenq= leaf.blocks[leaf.head].sumenq
212:     newBlock.sumdeq= leaf.blocks[leaf.head].sumdeq+1
213:     leaf.Append(newBlock)
214:     <b, i>= IndexDequeue(leaf.head, 1)
215:     output= FindResponse(b, i)
216:     return output
217: end Dequeue

218: <int, int> FindResponse(int b, int i)                   ▷ Returns the response to  $D_{root,b,i}$ .
219:     if root.blocks[b-1].size + root.blocks[b].numenq - i < 0 then           ▷ Check if the queue is empty.
220:         return null
221:     else
222:         e= i - root.blocks[b-1].size + root.blocks[b-1].sumenq             ▷  $E_e(root)$  is the response.
223:         return root.GetEnqueue(root.DSearch(e, b))
224:     end if
225: end FindResponse
```

---

---

**Algorithm Root**

---

$\rightsquigarrow$  Precondition:  $\text{root.blocks}[\text{end}].\text{sum}_{\text{enq}} \geq e$

801:  $\langle \text{int}, \text{int} \rangle$  DSearch( $\text{int } e, \text{int } \text{end}$ )  $\triangleright$  Returns  $\langle b, i \rangle$  such that  $E_e(\text{root}) = E_i(\text{root}, b)$ .

802:      $\text{start} = \text{end} - 1$

803:     **while**  $\text{root.blocks}[\text{start}].\text{sum}_{\text{enq}} \geq e$  **do**

804:          $\text{start} = \max(\text{start} - (\text{end} - \text{start}), 0)$

805:     **end while**

806:      $b = \text{root.BinarySearch}(\text{sum}_{\text{enq}}, e, \text{start}, \text{end})$

807:      $i = e - \text{root.blocks}[b-1].\text{sum}_{\text{enq}}$

808:     **return**  $\langle b, i \rangle$

809: **end** DSearch

---

---

**Algorithm Leaf**

---

601: *void* Append(*block blk*)  $\triangleright$  Only called by the owner of the leaf.

602:      $\text{blocks}[\text{head}] = \text{blk}$

603:      $\text{head} += 1$

604:      $\text{parent.Propagate}()$

605: **end** Append

---

---

**Algorithm** *Node*

---

```
301: void Propagate()
302:   if not Refresh() then
303:     Refresh()
304:   end if
305:   if this is not root then
306:     parent.Propagate()
307:   end if
308: end Propagate

309: boolean Refresh()
310:   h= head
311:   for each dir in {left, right} do
312:     hdir= dir.head
313:     if dir.blocks[hdir] != null then
314:       dir.Advance(hdir)
315:     end if
316:   end for
317:   new= CreateBlock(h)
318:   if new.num==0 then return true
319: end if
320: result= blocks[h].CAS(null, new)
321: this.Advance(h)
322: return result
323: end Refresh

324: void Advance(int h)
325:   hp= parent.head
326:   blocks[h].super.CAS(null, hp)
327:   head.CAS(h, h+1)
328: end Advance

329: int BinarySearch(field f, int i, int start, int end)
330:   return min{j: blocks[j].f ≥ i}
331: end BinarySearch

332: Block CreateBlock(int i)
333:   block new= new(block)
334:   for each dir in {left, right} do
335:     indexlast= dir.head-1
336:     indexprev= blocks[i-1].enddir
337:     new.enddir= indexlast
338:     blocklast= dir.blocks[indexlast]
339:     blockprev= dir.blocks[indexprev]
340:     new.sumenq-dir= blocks[i-1].sumenq-dir +
341:       blocklast.sumenq - blockprev.sumenq
341:     new.sumdeq-dir= blocks[i-1].sumdeq-dir +
342:       blocklast.sumdeq - blockprev.sumdeq
342:   end for
343:   if this is root then
344:     new.size = max(root.blocks[i-1].size + new.numenq
345:       - new.numdeq, 0)
345:   end if
346:   return new
347: end CreateBlock
```

---

**Algorithm Node**

---

$\rightsquigarrow$  Precondition:  $\text{blocks}[b].\text{num}_{\text{enq}} \geq i \geq 1$

```
401: element GetEnqueue(int b, int i) ▷ Returns the element of  $E_i(\text{this}, b)$ .
402:   if this is leaf then
403:     return blocks[b].element
404:   else if  $i \leq \text{blocks}[b].\text{num}_{\text{enq-left}}$  then ▷  $E_i(\text{this}, b)$  is in the left child of this node.
405:     subBlock= left.BinarySearch(sumenq, i+blocks[b-1].sumenq-left, blocks[b-1].endleft+1, blocks[b].endleft)
406:     return left.GetEnqueue(subBlock, i)
407:   else
408:     i= i-blocks[b].numenq-left
409:     subBlock= right.BinarySearch(sumenq, i+right.blocks[b-1].sumenq-right, blocks[b-1].endright+1, blocks[b].endright)
410:     return right.GetEnqueue(subBlock, i)
411:   end if
412: end GetEnqueue
```

$\rightsquigarrow$  Precondition: bth block of the node has propagated up to the root and  $\text{blocks}[b].\text{num}_{\text{deq}} \geq i$ .

```
413: <int, int> IndexDequeue(int b, int i) ▷ Returns <x, y> if  $D_i(\text{this}, b) = D_y(\text{root}, x)$ .
414:   if this is root then
415:     return <b, i>
416:   else
417:     dir= (parent.left==n ? left: right)
418:     sb= (parent.blocks[blocks[b].super].sumdeq-dir > blocks[b].sumdeq ? blocks[b].super: blocks[b].super+1)
419:     if dir is left then
420:       i+= blocks[b-1].sumdeq-parent.blocks[sb-1].sumdeq-left
421:     else
422:       i+= blocks[b-1].sumdeq-parent.blocks[sb-1].sumdeq-right
423:       i+= parent.blocks[sb].numdeq-left
424:     end if
425:     return this.parent.IndexDequeue(sb, i)
426:   end if
427: end IndexDequeue
```

---

## 4 Proof of Correctness

We adopt linearizability as our definition of correctness. In our case, where we create the linearization ordering in the root, we need to prove (1) the ordering is legal, i.e, for every execution on our queue if operation  $op_1$  terminates before operation  $op_2$  then  $op_1$  is linearized before operation  $op_2$  and (2) if we do operations sequentially in their the linearization order, operations get the same results as in our queue. The proof is structured like this. First, we define and prove some facts about blocks and the node's `head` field. Then, we introduce the linearization ordering formally. Next, we prove double **Refresh** on a node is enough to propagate its children's new operations up to the node, which is used to prove (1). After this, we prove some claims about the size and operations of each block, which we use to prove the correctness of `DSearch()`, `GetEnqueue()` and `IndexDequeue()`. Finally, we prove the correctness of the way we compute the response of a dequeue, which establishes (2).

### 4.1 Basic Properties

In this subsection we talk about some properties of blocks and fields of the tree nodes.

A block is an object storing some statistics, as described in Algorithm Queue. A block in a node implicitly represents a set of operations.

**Definition 1** (Ordering of a block in a node). Let  $b$  be  $n.\text{blocks}[i]$  and  $b'$  be  $n.\text{blocks}[j]$ . We call  $i$  the *index* of block  $b$ . Block  $b$  is *before* block  $b'$  in node  $n$  if and only if  $i < j$ . We define *the prefix* for block  $b$  in node  $n$  to be the blocks in  $n.\text{blocks}[0..i]$ .

Next, we show that the value of `head` in a node can only be increased. By the termination of a **Refresh**, `head` has been incremented by the process doing the **Refresh** or by another process.

**Observation 2.** *For each node  $n$ ,  $n.\text{head}$  is non-decreasing over time.*

*Proof.* The claim follows trivially from the code since `head` is only changed by incrementing in Line 327 of **Advance**. □

**Lemma 3.** *Let  $R$  be an instance of **Refresh** on a node  $n$ . After  $R$  terminates,  $n.\text{head}$  is greater than the value read in line 310 of  $R$ .*

*Proof.* If the **CAS** in Line 327 is successful then the claim holds. Otherwise  $n.\text{head}$  has changed from the value that was read in Line 310. By Observation 2 this means another process has incremented  $n.\text{head}$ . □

Now we show  $n.\text{blocks}[n.\text{head}]$  is either the last block written into node  $n$  or the first empty block in  $n$ .

**Invariant 4** (headPosition). If the value of  $n.\text{head}$  is  $h$  then  $n.\text{blocks}[i] = \text{null}$  for  $i > h$  and  $n.\text{blocks}[i] \neq \text{null}$  for  $0 \leq i < h$ .

*Proof.* Initially the invariant is true since  $n.\text{head} = 1$ ,  $n.\text{blocks}[0] \neq \text{null}$  and  $n.\text{blocks}[x] = \text{null}$  for every  $x > 0$ . The truth of the invariant may be affected by writing into  $n.\text{blocks}$  or incrementing  $n.\text{head}$ . We show that if the invariant holds before such a change then it still holds after the change.

In the algorithm,  $n.\text{blocks}$  is modified only on Line 320, which updates  $n.\text{blocks}[h]$  where  $h$  is the value read from  $n.\text{head}$  in Line 310. Since the CAS in Line 320 is successful it means  $n.\text{head}$  has not changed from  $h$  before doing the CAS: if  $n.\text{head}$  had changed before the CAS then it would be greater than  $h$  by Observation 2 and hence  $n.\text{blocks}[h] \neq \text{null}$  and by the induction hypothesis, so the CAS would fail. Writing into  $n.\text{blocks}[h]$  when  $h = n.\text{head}$  preserves the invariant, since the claim does not talk about the content of  $n.\text{blocks}[n.\text{head}]$ .

The value of  $n.\text{head}$  is modified only in Line 327 of **Advance**. If  $n.\text{head}$  is incremented to  $h + 1$  it is sufficient to show  $n.\text{blocks}[h] \neq \text{null}$ . **Advance** is called in Lines 314 and 321. For Line 314,  $n.\text{blocks}[h] \neq \text{null}$  because of the **if** condition in Line 313. For Line 321, Line 320 was finished before doing 321. Whether Line 320 is successful or not,  $n.\text{blocks}[h] \neq \text{null}$  after the  $n.\text{blocks}[h].\text{CAS}$ .  $\square$

We define the subblocks of a block recursively.

**Definition 5** (Subblock). A block is a *direct subblock* of the  $i$ th block in node  $n$  if it is in

$$n.\text{left}.\text{blocks}[n.\text{blocks}[i-1].\text{end}_{\text{left}}+1 \dots n.\text{blocks}[i].\text{end}_{\text{left}}]$$

or in

$$n.\text{right}.\text{blocks}[n.\text{blocks}[i-1].\text{end}_{\text{right}}+1 \dots n.\text{blocks}[i].\text{end}_{\text{right}}].$$

Block  $b$  is a *subblock* of block  $c$  if  $b$  is a direct subblock of  $c$  or a subblock of a direct subblock of  $c$ . We say block  $b$  is *propagated* to node  $n$  if  $b$  is in  $n.\text{blocks}$  or is a subblock of a block in  $n.\text{blocks}$ .

The next lemma is used to prove the subblocks of two blocks in a node are disjoint.

**Lemma 6.** If  $n.\text{blocks}[i] \neq \text{null}$  and  $i > 0$  then  $n.\text{blocks}[i].\text{end}_{\text{left}} \geq n.\text{blocks}[i-1].\text{end}_{\text{left}}$  and  $n.\text{blocks}[i].\text{end}_{\text{right}} \geq n.\text{blocks}[i-1].\text{end}_{\text{right}}$ .



*Proof.* Consider the block  $b$  written into  $n.\text{blocks}[i]$  by CAS at Line 320. Block  $b$  is created by the `CreateBlock( $i$ )` called at Line 317. Prior to this call to `CreateBlock( $i$ )`,  $n.\text{head} = i$  at Line 310, so  $n.\text{blocks}[i - 1]$  is already a non-null value  $b'$  by Invariant 4. Thus, the `CreateBlock( $i - 1$ )` that created  $b'$  terminated before the `CreateBlock( $i$ )` that creates  $b$  is invoked. The value written into  $b.\text{end}_{\text{left}}$  at Line 337 of `CreateBlock( $i$ )` was one less than the value read at Line 335 of `CreateBlock( $i$ )`. Similarly, the value in  $n.\text{blocks}[i - 1].\text{end}_{\text{left}}$  was one less than the value read from  $n.\text{left.head}$  during the call to `CreateBlock( $i - 1$ )`. By Observation 2,  $n.\text{left.head}$  is non-decreasing, so  $b'.\text{end}_{\text{left}} \leq b.\text{end}_{\text{left}}$ . The proof for  $\text{end}_{\text{right}}$  is similar.  $\square$

**Lemma 7.** *Subblocks of any two blocks in node  $n$  do not overlap.*

*Proof.* We are going to prove the lemma by contradiction. Consider the lowest node  $n$  in the tree that violates the claim. Then subblocks of  $n.\text{blocks}[i]$  and  $n.\text{blocks}[j]$  overlap for some  $i < j$ . Since  $n$  is the lowest node in the tree violating the claim, direct subblocks of blocks of  $n.\text{blocks}[i]$  and  $n.\text{blocks}[j]$  have to overlap. Without loss of generality assume left child subblocks of  $n.\text{blocks}[i]$  overlap with the left child subblocks of  $n.\text{blocks}[j]$ . By Lemma 6 we have  $n.\text{blocks}[i].\text{end}_{\text{left}} \leq n.\text{blocks}[j - 1].\text{end}_{\text{left}}$ , so the ranges  $[n.\text{blocks}[i - 1].\text{end}_{\text{left}} + 1 \dots n.\text{blocks}[i].\text{end}_{\text{left}}]$  and  $[n.\text{blocks}[j - 1].\text{end}_{\text{left}} + 1 \dots n.\text{blocks}[j].\text{end}_{\text{left}}]$  cannot overlap. Therefore, direct subblocks of  $n.\text{blocks}[i]$  and  $n.\text{blocks}[j]$  cannot overlap.  $\square$

**Definition 8** (Superblock). Block  $b$  is *superblock* of block  $c$  if  $c$  is a direct subblock of  $b$ .

**Corollary 9.** *Every block has at most one superblock.*

*Proof.* A block having more than one superblock contradicts Lemma 7.  $\square$

Now we can define the operations of a block using the definition of subblocks.

**Definition 10** (Operations of a block). A block  $b$  in a leaf represents an `Enqueue()` if  $b.\text{element} \neq \text{null}$ . Otherwise, if  $b.\text{element} = \text{null}$ ,  $b$  represents a `Dequeue()`. The set of operations of block  $b$  is the union of the operations in leaf subblocks of  $b$ . We denote the set of operations of block  $b$  by  $\text{ops}(b)$  and the union of operations of a set of blocks  $B$  by  $\text{ops}(B)$ . We also say  $b$  contains  $op$  if  $op \in \text{ops}(b)$ .

Operations are distinct `Enqueues` and `Dequeues` invoked by processes. The next lemma proves that each operation appears at most once in the blocks of a node.

**Lemma 11.** *If  $op$  is in  $n.\text{blocks}[i]$  then there is no  $j \neq i$  such that  $op$  is in  $n.\text{blocks}[j]$ .*

*Proof.* We prove this claim using Lemma 7. Assume  $op$  is in the subblocks of both  $n.blocks[i]$  and  $n.blocks[j]$ . From Corollary 7 we know that the subblocks of these blocks are different, so there are two leaf blocks containing  $op$ . Since each process puts each operation in only one block of its leaf then  $op$  cannot be in two leaf blocks. This is a contradiction.  $\square$

**Definition 12.**  $n.blocks[i]$  is *established* at time  $t$  if  $n.head > i$ . An operation is *established* in node  $n$  if it is in an established block of  $n$ .  $EST_t^n$  is the set of established operations in node  $n$  at time  $t$ .

Now we want to say the blocks of a node grow over time.

**Observation 13.** If time  $t < \text{time } t'$  ( $t$  is before  $t'$ ), then  $ops(n.blocks)$  at time  $t$  is a subset of  $ops(n.blocks)$  at time  $t'$ .

*Proof.* Blocks are only appended (not modified) with CAS to  $n.blocks[n.head]$ , so the set of blocks of a node after the CAS contains the the set of blocks before the CAS.  $\square$

**Corollary 14.** If time  $t < \text{time } t'$ , then  $EST_n^t \subseteq EST_n^{t'}$ .

*Proof.* From Observations 2, 13.  $\square$

## 4.2 Ordering Operations

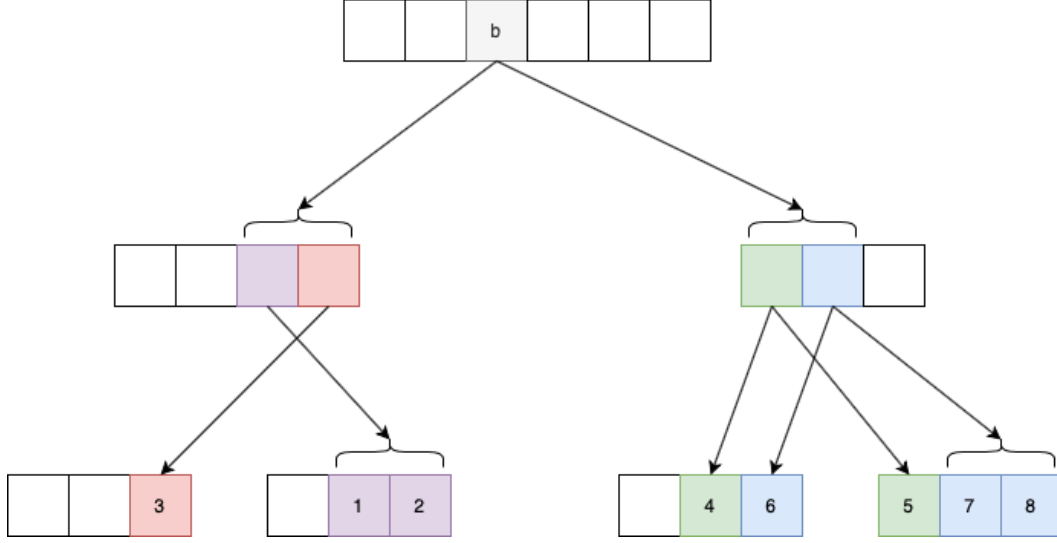


Figure 15: Order of operations in  $b$ . Operations in the leaves are ordered with numerical order shown in the drawing.

Now we define the ordering of operations stored in each node. In the non-root nodes we only need to order operations of a type among themselves. Processes are numbered from 1 to  $p$  and leaves of the tree are assigned from left to right. We will show in Lemma 28 that there is at most one operation from each process in a given block.

**Definition 15** (Ordering of operations inside the nodes).

- $E(n, b)$  is the sequence of enqueue operations in  $ops(n.blocks[b])$  defined recursively as follows.  
 $E(leaf, b)$  is the single enqueue operation in  $ops(leaf.blocks[b])$  or an empty sequence if  $leaf.blocks[b]$  represents a dequeue operation. If  $n$  is an internal node, then

$$E(n, b) = E(n.left, n.blocks[b-1].end_{left} + 1) \cdots E(n.left, n.blocks[b].end_{left}) \cdot \\ E(n.right, n.blocks[b-1].end_{right} + 1) \cdots E(n.right, n.blocks[b].end_{right}).$$

- $E_i(n, b)$  is the  $i$ th enqueue in  $E(n, b)$ .
- The order of the enqueue operations in the node  $n$  is  $E(n) = E(n, 1) \cdot E(n, 2) \cdot E(n, 3) \cdots$
- $E_i(n)$  is the  $i$ th enqueue in  $E(n)$ .

- $D(n, b)$  is the sequence of dequeue operations in  $ops(n.blocks[b])$  defined recursively as follows.  
 $D(leaf, b)$  is the single dequeue operation in  $ops(leaf.blocks[b])$  or an empty sequence if  $leaf.blocks[b]$  represents an enqueue operation. If  $n$  is an internal node, then

$$D(n, b) = D(n.left, n.blocks[b-1].end_{left} + 1) \cdots D(n.left, n.blocks[b].end_{left}) \cdot \\ D(n.right, n.blocks[b-1].end_{right} + 1) \cdots D(n.right, n.blocks[b].end_{right}).$$

- $D_i(n, b)$  is the  $i$ th enqueue in  $D(n, b)$ .
- The order of the dequeue operations in the node  $n$  is  $D(n) = D(n, 1) \cdot D(n, 2) \cdot D(n, 3) \dots$
- $D_i(n)$  is the  $i$ th dequeue in  $D(n)$ .

The linearization ordering is given by the order that operations appear in the blocks in the root.

**Definition 16** (Linearization).

$$L = E(root, 1) \cdot D(root, 1) \cdot E(root, 2) \cdot D(root, 2) \cdot E(root, 3) \cdot D(root, 3) \cdots$$

**Observation 17.** For any node  $n$  and indices  $i < j$  of blocks in  $n$ , we have

$$n.blocks[j].sum_x - n.blocks[i].sum_x = \sum_{k=i+1}^j n.blocks[k].num_x$$

where  $x$  in  $\{enq, deq, enq-left, enq-right, deq-left, deq-right\}$ .

Next claim is also true if we replace **enq** with **deq** and  $E$  with  $D$ .

**Lemma 18.** Let  $B, B'$  be  $n.blocks[b], n.blocks[b-1]$  respectively.

- (1) If  $n$  is an internal node  $B.num_{enq-left} = \left| E(n.left, B'.end_{left} + 1) \cdots E(n.left, B.end_{left}) \right|$ .
- (2) If  $n$  is an internal node  $B.num_{enq-right} = \left| E(n.right, B'.end_{right} + 1) \cdots E(n.right, B.end_{right}) \right|$ .
- (3)  $B.num_{enq} = \left| E(n, b) \right|$ .

*Proof.* We prove the claim by induction on height of node  $n$ . Base case (3) for leaves is trivial. Supposing

the claim is true for  $n$ 's children, we prove the correctness of the claim for  $n$ .

$$\begin{aligned}
B.\text{num}_{\text{enq-left}} &= B.\text{sum}_{\text{enq-left}} - B'.\text{sum}_{\text{enq-left}} && \text{Definition of } \text{num}_{\text{enq}} \\
&= B'.\text{sum}_{\text{enq-left}} + n.\text{left.blocks}[B.\text{end}_{\text{left}}].\text{sum}_{\text{enq}} \\
&\quad - n.\text{left.blocks}[B'.\text{end}_{\text{left}}].\text{sum}_{\text{enq}} - B'.\text{sum}_{\text{enq-left}} && \text{CreateBlock} \\
&= n.\text{left.blocks}[B.\text{end}_{\text{left}}].\text{sum}_{\text{enq}} - n.\text{left.blocks}[B'.\text{end}_{\text{left}}].\text{sum}_{\text{enq}} \\
&= \sum_{i=B'.\text{end}_{\text{left}}+1}^{B.\text{end}_{\text{left}}} n.\text{left.blocks}[i].\text{num}_{\text{enq}} && \text{Observation 17} \\
&= \left| E(n.\text{left}, B'.\text{end}_{\text{left}} + 1) \cdots E(n.\text{left}, B.\text{end}_{\text{left}}) \right| && \text{Induction hypothesis (3)}
\end{aligned}$$

The last line holds because of the induction hypothesis (3). (2) is similar to (1). Now we prove (3) starting from the Definition of  $E(n, b)$ .

$$\begin{aligned}
E(n, b) &= E(n.\text{left}, n.\text{blocks}[b-1].\text{end}_{\text{left}} + 1) \cdots E(n.\text{left}, n.\text{blocks}[b].\text{end}_{\text{left}}) \cdot \\
&\quad E(n.\text{right}, n.\text{blocks}[b-1].\text{end}_{\text{right}} + 1) \cdots E(n.\text{right}, n.\text{blocks}[b].\text{end}_{\text{right}}).
\end{aligned}$$

By (1) and (2) we have  $\left| E(n, b) \right| = B.\text{num}_{\text{enq-left}} + B.\text{num}_{\text{enq-right}} = B.\text{num}_{\text{enq}}$ . □

Next claim is also true if we replace **enq** with **deq** and  $E$  with  $D$ .

**Corollary 19.** *Let  $B$  be  $n.\text{blocks}[b]$  and **enq** be in  $\{\text{enq}, \text{deq}\}$ .*

- (1) *If  $n$  is an internal node  $B.\text{sum}_{\text{enq-left}} = \left| E(n.\text{left}, 1) \cdots E(n.\text{left}, B.\text{end}_{\text{left}}) \right|$*
- (2) *If  $n$  is an internal node  $B.\text{sum}_{\text{enq-right}} = \left| E(n.\text{right}, 1) \cdots E(n.\text{right}, B.\text{end}_{\text{right}}) \right|$*
- (3)  *$B.\text{sum}_{\text{enq}} = \left| E(n, 1) \cdot E(n, 2) \cdots E(n, b) \right|$*

### 4.3 Propagating Operations to the Root

In this section we explain why two **Refreshes** are enough to propagate a nodes operations to its parent.

**Definition 20.** Let  $t^{op}$  be the time  $op$  is invoked,  $^{op}t$  be the time  $op$  terminates,  $t_l^{op}$  be the time immediately before running Line  $l$  of operation  $op$  and  $^{op}_l t$  be the time immediately after running Line  $l$  of operation  $op$ . We sometimes suppress  $op$  and write  $t_l$  or  $_l t$  if  $op$  is clear in the context. In the text  $v_l$  is the value of variable  $v$  immediately after line  $l$  for the process we are talking about and  $v_t$  is the value of variable  $v$  at time  $t$ .

**Definition 21** (Successful Refresh). An instance of **Refresh** is *successful* if its **CAS** in Line 320 returns **true**. If a successful instance of **Refresh** terminates, we say it is *complete*.

In the next two results we show for every successful **Refresh**, all the operations established in the children before the **Refresh** are in the parent after the **Refresh**'s successful **CAS** at Line 320.

**Lemma 22.** *If  $R$  is a successful instance of  $n$ .Refresh, then we have  $EST_{n.\text{left}}^{t^R} \cup EST_{n.\text{right}}^{t^R} \subseteq ops(n.\text{blocks}_{320})$ .*

*Proof.* We show

$$\begin{aligned} EST_{n.\text{left}}^{t^R} &= ops(n.\text{left.blocks}[0..n.\text{left.head}_{309} - 1]) \\ &\subseteq ops(n.\text{blocks}_{320}) = ops(n.\text{blocks}[0..n.\text{head}_{320}]). \end{aligned}$$

Line 320 stores a block **new** in  $n$  that has  $\text{end}_{\text{left}} = n.\text{left.head}_{335} - 1$ . Therefore, by Definition 5, after the successful **CAS** in Line 320 we know all blocks in  $n.\text{left.blocks}[1 \dots n.\text{left.head}_{335} - 1]$  are subblocks of  $n.\text{blocks}[1 \dots n.\text{head}_{310}]$ . Because of Lemma 2 we have  $n.\text{left.head}_{309} - 1 < n.\text{left.head}_{335} - 1$  and  $n.\text{head}_{310} < n.\text{head}_{320}$ . From Observation 13 the claim follows. The proof for the right child is the same.  $\square$

**Corollary 23.** *If  $R$  is a complete instance  $n$ .Refresh, then we have  $EST_{n.\text{left}}^{t^R} \cup EST_{n.\text{right}}^{t^R} \subseteq EST_n^{Rt}$ .*

*Proof.* The left hand side is the same as Lemma 22, so it is sufficient to show when  $R$  terminates the established blocks in  $n$  are a superset of  $n.\text{blocks}_{320}$ . Line 320 writes the block **new** in  $n.\text{blocks}[h]$  where  $h$  is value of  $n.\text{head}$  read at Line 310. Because of Lemma 3 we are sure that  $n.\text{head} > h$  when  $R$  terminates. So the block **new** appended to  $n$  at Line 320 is established at  $^R t$ .  $\square$

In the next lemma we show that if two consecutive instances of **Refresh** by the same process on node  $n$  fail, then the blocks established in the children of  $n$  before the first **Refresh** are guaranteed to be in  $n$  after the second **Refresh**.

**Lemma 24.** *Consider two consecutive terminating instances  $R_1, R_2$  of **Refresh** on internal node  $n$  by process  $p$ . If neither  $R_1$  nor  $R_2$  is a successful **Refresh**, then we have  $EST_{n.\text{left}}^{t^{R_1}} \cup EST_{n.\text{right}}^{t^{R_1}} \subseteq EST_n^{R_2 t}$ .*

*Proof.* Let  $R_1$  read  $i$  from  $n.\text{head}$  at Line 310. By Lemma 3,  $R_1$  and  $R_2$  both cannot read the same value  $i$ . By Observation 2,  $R_2$  reads a larger value of  $n.\text{head}$  than  $R_1$ .

Consider the case where  $R_1$  reads  $i$  and  $R_2$  reads  $i + 1$  from Line 310. As  $R_2$ 's CAS in Line 320 returns **false**, there is another successful instance  $R'_2$  of  $n.\text{Refresh}$  that has done a CAS successfully into  $n.\text{blocks}[i + 1]$  before  $R_2$  tries to CAS.  $R'_2$  creates its block **new** after reading the value  $i + 1$  from  $n.\text{head}$  (Line 310) and  $R_1$  reads the value  $i$  from  $n.\text{head}$ . By Observation 2 we have  $R_1 t < t_{310}^{R_1} < t_{310}^{R'_2}$  (see Figure 16). By Lemma 23 we have  $EST_{R'_2 t}^{n.\text{left}} \cup EST_{R'_2 t}^{n.\text{right}} \subseteq ops(n.\text{blocks}_{R'_2 t})$ . Also by Lemma 3 on  $R_2$ , the value of  $n.\text{head}$  is more than  $i + 1$  after  $R_2$  terminates, so the block appended by  $R'_2$  to  $n$  is established by the time  $R_2$  terminates. To summarize,  $R_1 t$  is before  $R'_2$ 's read of  $n.\text{head}$  ( $t_{310}^{R'_2}$ ) and  $R'_2$ 's successful CAS ( $t_{320}^{R'_2}$ ) is before  $R_2$ 's termination ( $t^{R_2}$ ), so by Observation and Lemma 3 we have  $EST_{n.\text{left}}^{t^{R_1}} \cup EST_{n.\text{right}}^{t^{R_1}} \subseteq ops(n.\text{blocks}_{t_{320}^{R'_2}}) \subseteq EST_{n.\text{left}}^{R_2 t}$ .

If  $R_2$  reads some value greater than  $i + 1$  in Line 310 it means  $n.\text{head}$  has been incremented more than two times since  $R_1 t$ . By Lemma 4, when  $n.\text{head}$  is incremented from  $i + 1$  to  $i + 2$ ,  $n.\text{blocks}[i + 1]$  is non-null. Let  $R_3$  be the **Refresh** on  $n$  that has put the block in  $n.\text{blocks}[i + 1]$ .  $R_3$  read  $n.\text{head} = i + 1$  at Line 310 and has put its block in  $n.\text{blocks}[i + 1]$  before  $R_2$ 's read of  $n.\text{head}$  at Line 310. So we have  $t^{R_1} <_{310}^{R_3} t <_{320}^{R_3} t < t_{310}^{R_2} <_2^R t$ . From Observation 13 on the operations before and after  $R_3$ 's CAS and Lemmas 22, 3 on  $R_3$  the claim holds.  $\square$

**Corollary 25.**  $EST_{n.\text{left}}^{302t} \cup EST_{n.\text{right}}^{302t} \subseteq EST_n^{t_{303}}$

*Proof.* If the first **Refresh** in line 302 returns **true** then by Lemma 23 the claim holds. If the first **Refresh** failed and the second **Refresh** succeeded the claim still holds by Lemma 23. Otherwise both failed and the claim is satisfied by Lemma 24.  $\square$

Now we show that after **Append**( $b$ ) on a leaf finishes, the operation contained in  $b$  will be established in **root**.

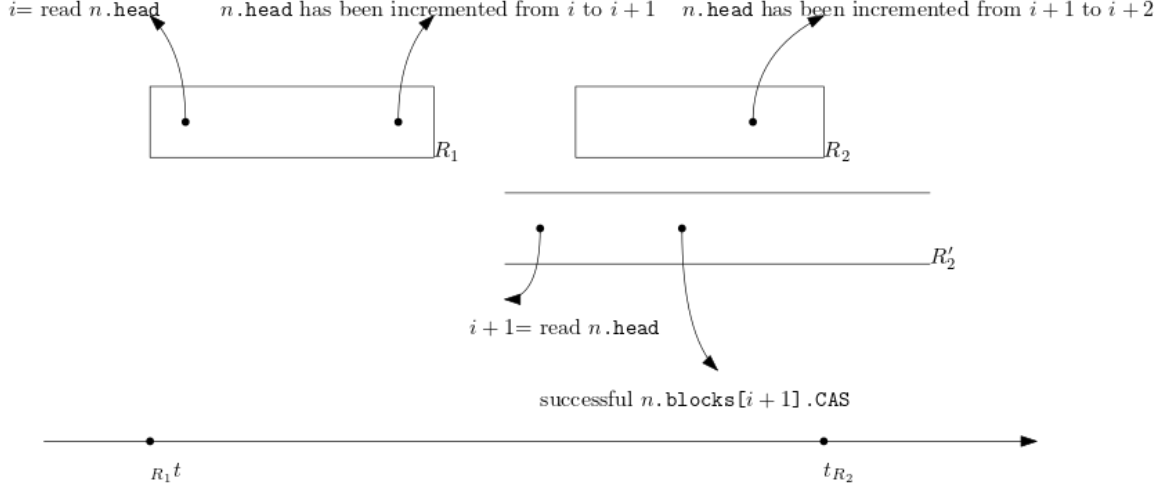


Figure 16:  $R_1t < t_{310}^{R_1} < \text{incrementing } n.\text{head} \text{ from } i \text{ to } i + 1 < t_{310}^{R'_2} < t_{320}^{R'_2} < \text{incrementing } n.\text{head} \text{ from } i + 1 \text{ to } i + 2 < t_{R_2}$

**Corollary 26.** For  $A = l.\text{Append}(b)$  we have  $\text{ops}(b) \subseteq \text{EST}_n^{t^A}$  for each node  $n$  in the path from  $l$  to root.

*Proof.*  $A$  adds  $b$  to the assigned leaf of the process, establishes it at Line 603 and then calls **Propagate** on the parent of the leaf where it appended  $b$ . For every node  $n$ ,  $n.\text{Propagate}$  appends  $b$  to  $n$ , establishes it in  $n$  by Corollary 25 and then calls  $n.\text{parent}.\text{Propagate}$  until  $n$  is root.  $\square$

**Corollary 27.** After  $l.\text{Append}(b)$  finishes,  $b$  is subblock of exactly one block in each node along the path from  $l$  to the root.

*Proof.* By the previous corollary and Lemma 27 there is exactly one block in each node containing  $b$ .  $\square$



## 4.4 Correctness of GetEnqueue

First we prove some claims about the size and operations of a block. These lemmas will be used later for the correctness and analysis of `GetEnqueue()`.

**Lemma 28.** *Each block contains at most one operation of each process, and therefore at most  $p$  operations in total.*

*Proof.* To derive a contradiction, assume there are two operations  $op_1$  and  $op_2$  of process  $p$  in block  $b$  in node  $n$ . Without loss of generality  $op_1$  is invoked earlier than  $op_2$ . Process  $p$  cannot invoke more than one operation concurrently, so  $op_1$  has to be finished before  $op_2$  begins. By Corollary 27, before  $op_2$  calls `Append`,  $op_1$  exists in every node of the tree on the path from  $p$ 's leaf to the root. Since  $b$  contains  $op_2$ , it must be created after  $op_2$  is invoked. This means there is some block  $b'$  before  $b$  in  $n$  containing  $op_1$ . The existence of  $op_1$  in  $b$  and  $b'$  contradicts Lemma 11.  $\square$

**Lemma 29.** *Each block has at most  $p$  direct subblocks.*

*Proof.* The claim follows directly from Lemma 28 and the observation that each block appended to an internal node contains at least one operation, due to the test on Line 318. We can also see the blocks in the leaves have exactly one operation in the `Enqueue()` and `Dequeue()` routines.  $\square$

`DSearch( $e$ ,  $end$ )` returns a pair  $\langle b, i \rangle$  such that the  $i$ th `Enqueue` in the  $b$ th block of the root is the  $e$ th `Enqueue` in the entire sequence stored in the root.

**Lemma 30** (DSearch Correctness). *If  $\text{root.blocks}[end] \neq \text{null}$  and  $1 \leq e \leq \text{root.blocks}[end].\text{sum}_{\text{enq}}$ , `DSearch( $e$ ,  $end$ )` returns  $\langle b, i \rangle$  such that  $E_i(\text{root}, b) = E_e(\text{root})$ .*

*Proof.* From Lines 340 and 341 we know the `sumenq-left` and `sumenq-right` fields of `blocks` in each node are sorted in non-decreasing order. Since `sumenq = sumenq-left + sumenq-right`, the `sumenq` values of `root.blocks[0..end]` are also non-decreasing. Furthermore, since `root.blocks[0].sumenq = 0` and `root.blocks[end].sumenq ≥  $e$` , there is a  $b$  such that `root.blocks[b].sumenq ≥  $e$`  and `root.blocks[b-1].sumenq <  $e$`  by Lemma 19. Block `root.blocks[b]` contains  $E_i(\text{root}, b)$ . Lines 802–805 doubles the search range in Line 804 and will eventually reach `start` such that `root.blocks[start].sumenq ≤  $e$  ≤ root.blocks[end].sumenq`. Then, in Line 806, the binary search finds the  $b$  such that `root.blocks[b-1].sumenq <  $e$  ≤ root.blocks[b].sumenq`. By Corollary 19, `root.blocks[b]` is the block that contains  $E_e(\text{root})$ . Finally  $i$  is computed using the definition of `sumenq` and Corollary 19.  $\square$

**Lemma 31** (GetEnqueue correctness). *If  $1 \leq i \leq n.\text{blocks}[b].\text{num}_{\text{enq}}$  then  $n.\text{GetEnqueue}(b, i)$  returns  $E_i(n, b).\text{element}$ .*

*Proof.* We are going to prove this lemma by induction on the height of node  $n$ . For the base case, suppose  $n$  is a leaf. Leaf blocks each contain exactly one operation,  $n.\text{blocks}[b].\text{sum}_{\text{enq}} \leq 1$ , which means only  $n.\text{GetEnqueue}(b, 1)$  can be called when  $n$  is a leaf. Line 403 of  $n.\text{GetEnqueue}(b, 1)$  returns the `element` of the `Enqueue` operation stored in the  $b$ th block of leaf  $n$ , as required.

For the induction step we prove if  $n.\text{child}.\text{GetEnqueue}(b', i)$  returns  $E_i(n.\text{child}, b')$  then  $n.\text{GetEnqueue}(b, i)$  returns  $E_i(n, b)$ . From Definition 15 of  $E(n, b)$ , so operations from the left subblocks come before the operations from the right subblocks in a block (see Figure 17). By Observation 18, the  $\text{num}_{\text{enq-left}}$  field in  $n.\text{blocks}[b]$  is the number of `Enqueue()` operations from the blocks's subblocks in the left child of  $n$ . So the  $i$ th `Enqueue` operation in  $n.\text{blocks}[b]$  is propagated from the right child if and only if  $i$  is greater than  $n.\text{blocks}[b].\text{num}_{\text{enq-left}}$ . Line 404 decides whether the  $i$ th enqueue in the  $b$ th block of internal node  $n$  is in the left child or right child subblocks of  $n.\text{blocks}[b]$ . By Definitions 5 and 10 to find an operation in the subblocks of  $n.\text{blocks}[i]$  we need to search in the range

$$\begin{aligned} & n.\text{left}.\text{blocks}[n.\text{blocks}[i-1].\text{end}_{\text{left}}+1..n.\text{blocks}[i].\text{end}_{\text{left}}] \text{ or} \\ & n.\text{right}.\text{blocks}[n.\text{blocks}[i-1].\text{end}_{\text{right}}+1..n.\text{blocks}[i].\text{end}_{\text{right}}]. \end{aligned}$$

First we consider the case where the `Enqueue` we are looking for is in the left child. There are  $eb = n.\text{blocks}[b-1].\text{sum}_{\text{enq-left}}$  `Enqueues` in the blocks of  $n.\text{left}$  before the left subblocks of  $n.\text{blocks}[b]$ , so  $E_i(n, b)$  is  $E_{i+eb}(n.\text{left})$  which is  $E_{i'}(n.\text{left}, b')$  for some  $b'$  and  $i'$ . We can compute  $b'$  and then search for the  $i + eb$ th enqueue in  $n.\text{left}$ , where  $i'$  is  $i + eb - n.\text{left}.\text{blocks}[b'-1].\text{sum}_{\text{enq}}$ . The parameters in Line 405 are for searching  $E_{i+eb}(n.\text{left})$  in  $n.\text{left}.\text{blocks}$  in the range of left subblocks of  $n.\text{blocks}[b]$ , so this `BinarySearch` returns the index of the subblock containing  $E_i(n, b)$ .

Otherwise, the enqueue we are looking for is in the right child. Because `Enqueues` from the left subblocks are ordered before the `Enqueues` from the right subblocks, there are  $n.\text{blocks}[b].\text{num}_{\text{enq-left}}$  enqueues ahead of  $E_i(n, b)$  from the left child. So we need to search for  $i - n.\text{blocks}[b].\text{num}_{\text{enq-left}} + n.\text{blocks}[b-1].\text{sum}_{\text{enq-right}}$  in the right child (Line 409). Other parameters for the right child are chosen similarly to the left child.

So, in both cases the direct subblock containing  $E_i(n, b)$  is computed in Line 405 or 409. Finally,  $n.\text{child}.\text{GetEnqueue}(\text{subblock}, i)$  is invoked on the subblock containing  $E_i(n, b)$  and it returns  $E_i(n, b).\text{element}$

by the hypothesis of the induction. □

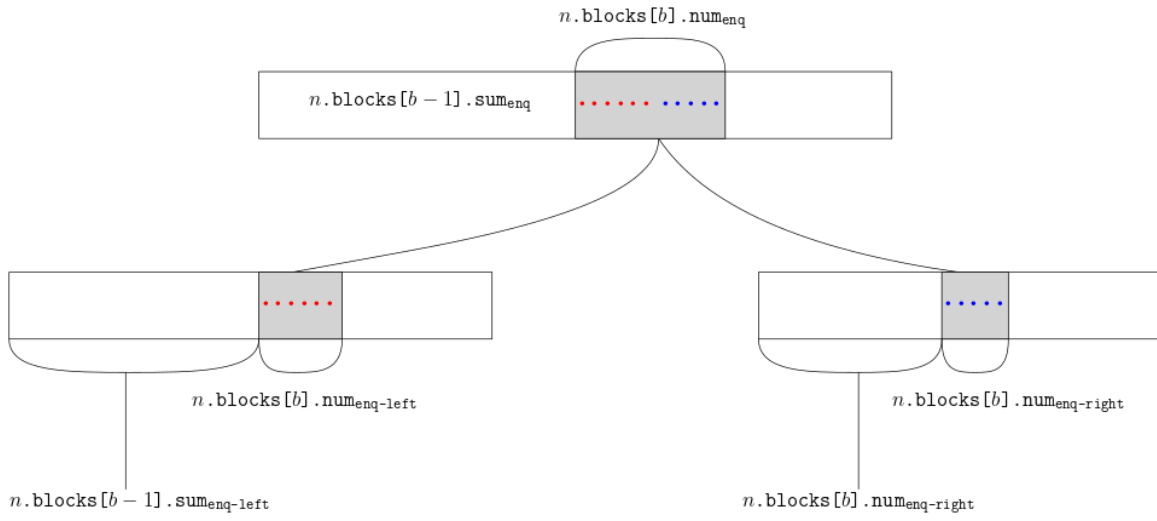


Figure 17: The number and ordering of the enqueue operations propagated from the left and the right child to  $n.\text{blocks}[b]$ . Both  $n.\text{blocks}[b]$  and its subblocks are shown in grey. **Enqueue** operations from the left subblocks (colored red), are ordered before the **Enqueue** operations from the right child (colored blue).

## 4.5 Correctness of IndexDequeue

The next few results show that the `super` field of a block is accurate within one of the actual index of the block's superblock in the parent node. Then we explain how it is used to compute the rank of a given dequeue in the root.

**Definition 32.** If a `Refresh` instance  $R_1$  does its CAS at Line 320 earlier than `Refresh` instance  $R_2$  we say  $R_1$  has *happened before*  $R_2$ .

**Observation 33.** After  $n.\text{blocks}[i].\text{CAS}(\text{null}, B)$  succeeds,  $n.\text{head}$  cannot increase from  $i$  to  $i + 1$  until  $B.\text{super}$  is set.

*Proof.* From Observation 2 we know the  $n.\text{head}$  changes only by the increment on Line 327. Before an instance of `Advance` increments  $n.\text{head}$  on Line 327, Line 326 ensures that  $n.\text{blocks}[\text{head}].\text{super}$  was set at Line 326. □

**Corollary 34.** If  $n.\text{blocks}[i].\text{super}$  is null, then  $n.\text{head} \leq i$  and  $n.\text{blocks}[i + 1]$  is null.

*Proof.* By Lemma 4 and Observation 33. □

Now let us consider how the `Refreshes` that took place on the parent of node  $n$  after block  $B$  was stored in  $n$  will help to set  $B.\text{super}$  and propagate  $B$  to the parent.

**Observation 35.** If the block created by an instance  $R_p$  of  $n.\text{parent.Refresh}$  contains block  $B = n.\text{blocks}[b]$  then  $R_p$  reads a value greater than  $b$  from  $n.\text{head}$  in Line 335.

**Lemma 36.** If  $B = n.\text{blocks}[b]$  is a direct subblock of  $n.\text{parent.blocks}[sb]$  then  $B.\text{super} \leq sb$ .

*Proof.* Let  $R_p$  be the instance of  $n.\text{parent.Refresh}$  that stores  $n.\text{parent.blocks}[sb]$ . By 35 if  $R_p$  propagates  $B$  it has to read a greater value than  $b$  from  $n.\text{head}$ , which means  $n.\text{head}$  was incremented from  $b$  to  $b + 1$  in Line 327. By Observation 33  $B.\text{super}$  was already set in Line 326. The value written in  $B.\text{super}$  was read in Line 325, before the CAS that sets  $B.\text{super}$ . From Observation 2 we know  $n.\text{parent.head}$  is non-decreasing so  $B.\text{super} \leq sb$ , since  $n.\text{parent.head}$  is still equal to  $sb$  when  $R_p$  executes its CAS at Line 320 by Invariant 6. The reader may wonder when the case  $b.\text{super} = sb$  happens. This can happen when  $n.\text{parent.blocks}[B.\text{super}] = \text{null}$  when  $B.\text{super}$  is written and  $R_p$  puts its created block into  $n.\text{parent.blocks}[B.\text{super}]$  afterwards. □

**Lemma 37.** *Let  $R_n$  be a Refresh that puts  $B$  in  $n.\text{blocks}[b]$  at Line 320. Then, the block created by one of the next two successful  $n.\text{parent.Refresh}$ es according to Definition 32 contains  $B$  and  $B.\text{super}$  is set when the second successful  $n.\text{parent.Refresh}$  reaches Line 317.*

*Proof.* Let  $R_{p1}$  be the first successful  $n.\text{parent.Refresh}$  after  $R_n$  and  $R_{p2}$  be the second next successful  $n.\text{parent.Refresh}$ . To derive a contradiction assume  $B$  was not propagated to  $n.\text{parent}$  by  $R_{p1}$  nor by  $R_{p2}$ .

Since  $R_{p2}$ 's created block does not contain  $B$ , by Observation 35 the value  $R_{p2}$  reads from  $n.\text{head}$  in Line 335 is at most  $b$ . From Observation 2 the value  $R_{p2}$  reads in Line 312 is also at most  $b$ .

$R_n$  puts  $B$  into  $n.\text{blocks}[b]$  so  $R_n$  reads the value  $b$  from  $n.\text{head}$ . Since  $R_{p2}$ 's CAS into  $n.\text{parent.blocks}$  is successful there should be a Refresh instance  $R'_p$  on  $n.\text{parent}$  that increments  $n.\text{parent}$  (Line 327) after  $R_{p1}$ 's Line 320 and before  $R_{p2}$ 's Line 310. We assumed  $t_{320}^{R_n} < t_{320}^{R_{p1}} < t_{320}^{R_{p2}}$  by Definition 32. Finally, Line 312 is after Line 310 and  $R_{p2}$ 's 310 is after  $R'_p$ 's Line 327, which is after  $R_n$ 's  $n.\text{blocks.CAS}$ .

$$\left. \begin{array}{l} R_n t_{320} <_{R_{p1}} t \\ R_{p1} t_{320} <_{R_{p'}} t <_{R_{p2}} t \\ R_{p2} t_{310} <_{R_{p2}} t \end{array} \right\} \Rightarrow R_n t_{320} <_{R_{p2}} t$$

So  $R_{p2}$  reads a value greater than or equal to  $b$  for  $n.\text{head}$  by Lemma 2.

Therefore  $R_{p2}$  reads  $n.\text{head} = b$ .  $R_{p2}$  calls  $n.\text{Advance}$  at Line Line 314, which ensures  $n.\text{head}$  is incremented from  $b$ . So the value  $R_{p2}$  reads in Line 335 of **CreateBlock** is greater than  $b$  and  $R_{p2}$ 's created block contains  $B$ . This is contradiction with our hypothesis.

Furthermore, if  $B.\text{super}$  was not set earlier it is set by  $R_{p2}$  call to  $n.\text{Advance}$  invoked from Line 314. □

**Corollary 38.** *If  $B = n.\text{blocks}[b]$  is propagated to  $n.\text{parent}$ , then  $B.\text{super}$  is equal to or one less than the index of the superblock of  $B$ .*

*Proof.* Let  $R_n$  be the  $n.\text{Refresh}$  that put  $B$  in  $n.\text{blocks}$  and let  $R_{p1}$  be the first successful  $n.\text{parent.Refresh}$  after  $R_n$  and  $R_{p2}$  be the second next successful  $n.\text{parent.Refresh}$ . Before  $B$  can be propagated to  $n$ 's parent,  $n.\text{head}$  must be greater than  $b$ , so by Observation 33  $B.\text{super}$  is set. From thr previous Lemma we know that  $B$  is propagated by second next successful Refresh's CAS on  $n.\text{parent.blocks}$ . To summarize we have  $n.\text{parent.head}_{R_{p2} t_{320}} = n.\text{parent.head}_{R_{p1} t_{320}} + 1$  and by Definition 32 and Observation 2

$n.\text{parent}.\text{head}_{\frac{R_{p^1}}{320}t} \leq n.\text{parent}.\text{head}_{\frac{R_n}{320}t}$ . The value that is set in  $B.\text{super}$  is read from  $n.\text{parent}.\text{head}$  after  $\frac{R_n}{320}t$ . So  $B.\text{super}$  is equal to or one less than the index of the superblock of  $B$ .  $\square$

Now using Corollary 38 on each step of the `IndexDequeue` we prove its correctness.

**Lemma 39** (`IndexDequeue` correctness). *If  $1 \leq i \leq n.\text{blocks}[b].\text{num\_deq}$  then  $n.\text{IndexDequeue}(b, i)$  returns  $\langle x, y \rangle$  such that  $D_i(n, b) = D_y(\text{root}, x)$ .*

*Proof.* We will prove this by induction on the distance of  $n$  from the `root`. The base case where  $n$  is `root` is trivial (see Line 415). For the non-root nodes  $n.\text{IndexDequeue}(b, i)$  computes  $sb$ , the index of the superblock of the  $b$ th block in  $n$ , in Line 418 by Corollary 38. After that, the position of  $D_i(n, b)$  in  $D(n.\text{parent}, sb)$  is computed in Lines 419–424. By Definition 15, `Dequeues` in a block are ordered based on the order of its subblocks from left to right. If  $D_i(n, b)$  was propagated from the left child, the number of dequeues in the left subblocks of  $n.\text{parent}.\text{blocks}[sb]$  before  $n.\text{blocks}[b]$  is considered in Line 420 (see Figure 18). Otherwise, if  $D_i(n, b)$  was propagated from the right child, the number of dequeues in the subblocks from the left child is considered to be ahead of the computed index (Line 421) (see Figure 19). Finally `IndexDequeue` is called on  $n.\text{parent}$  recursively and it returns the correct response by induction hypothesis.  $\square$

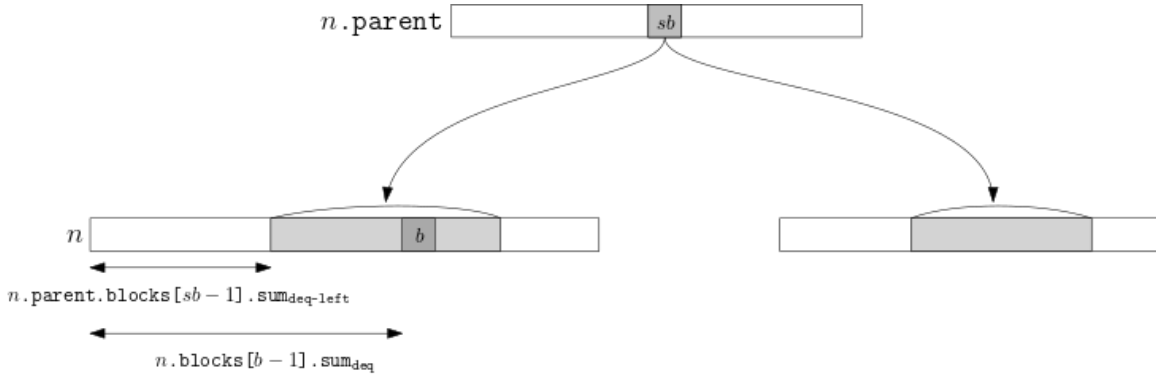


Figure 18: The number of `dequeue` operations before  $E_i(n, b)$  shown in the case where  $n$  is a left child.

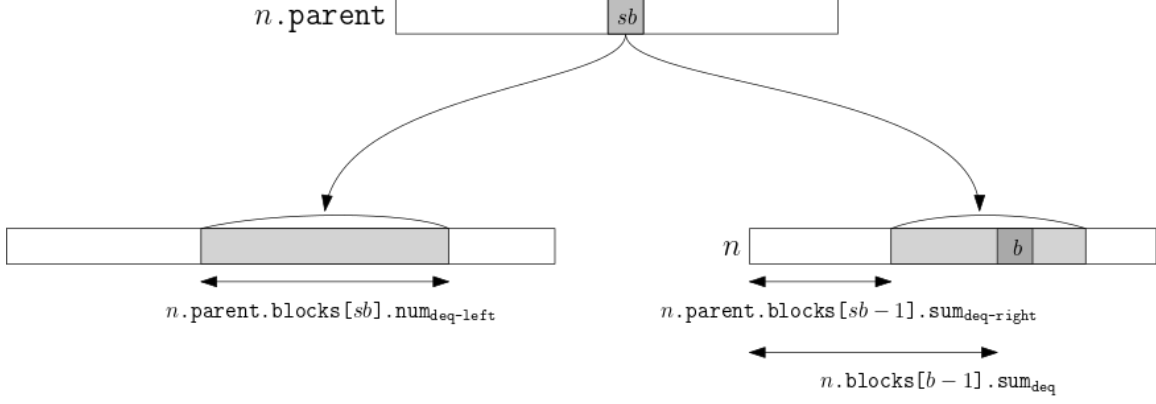


Figure 19: The number of **Dequeue** operations before  $E_i(n, b)$  shown in the case where  $n$  is a right child.

## 4.6 Linearizability

We now prove the two properties needed for linearizability.

**Lemma 40.**  *$L$  is a legal linearization ordering.*

*Proof.* We must show that, every operation that terminates is in  $L$  exactly once and if  $op_1$  terminates before  $op_2$  starts in execution then  $op_1$  is before  $op_2$  in the linearization. The first claim is directly reasoned from Lemma 27. For the latter, if  $op_1$  terminates before  $op_2$  starts,  $op_1$ .**Append** has terminated before  $op_2$ .**Append** started. From Lemma 26,  $op_1$  is in **root.blocks** before  $op_2$  starts to propagate. By definition of  $L$ ,  $op_1$  is linearized before  $op_2$ .  $\square$

Once some operations are aggregated in one block, they will get propagated up to the root together and they can be linearized in any order among themselves. We have chosen to put **Enqueues** in a block before **Edequeues** (see Definition 15).

**Definition 41.** If a **Dequeue** operation returns null it is called a *null Dequeue*, otherwise it is called *non-null Dequeue*.

Next we define the responses that **Dequeues** should return, according to the linearization.

**Definition 42.** Assume the operations in **root.blocks** are applied sequentially on an empty queue in the order of  $L$ .  $Resp(d) = e.element$  if the element of **Enqueue**  $e$  is the response to **Dequeue**  $d$ . Otherwise if  $d$  is a null dequeue then  $Resp(d) = null$ .

In the next lemma we show that the **size** field in each **root block** is computed correctly.

**Lemma 43.** `root.blocks[b].size` is the size of the queue if the operations in `root.blocks[0...b]` are applied in the order of  $L$ .

*Proof.* We prove the claim by induction on  $b$ . The base case when  $b = 0$  is trivial since the queue is initially empty and `root.blocks[0].size` = 0. We are going to show the correctness when  $b = i$  assuming correctness when  $b = i - 1$ . By Definition 15 **Enqueue** operations come before **Dequeue** operations in a block. By Lemma 18 `numenq` and `numdeq` fields in a block show the number of **Enqueue** and **Dequeue** operations in it. If there are more than `root.blocks[i - 1].size + root.blocks[i].numenq` dequeue operations in `root.blocks[i]` then the queue would become empty after `root.blocks[i]`. Otherwise the size of the queue after the  $b$ th block in the root is `root.blocks[b - 1].size + root.blocks[b].numenq - root.blocks[b].numdeq`. In both cases, this is same as the assignment on Line 344.  $\square$

The next lemma is useful to compute the number of non-null dequeues.

**Lemma 44.** If operations in the root are applied with the order of  $L$ , the number of non-null **Dequeues** in `root.blocks[0...b]` is `root.blocks[b].sumenq - root.blocks[b].size`.

*Proof.* There are `root.blocks[b].sumenq` enqueue operations in `root.blocks[0...b]`. The size of the queue after doing `root.blocks[0...b]` in order  $L$  is the number of *enqueues* in `root.blocks[0...b]` minus the number of *non-null Dequeues* in `root.blocks[0...b]`. By the correctness of the `size` field from Lemma 43 and `sumenq` field from Lemma 18, the number of *non-null Dequeues* is `root.blocks[b].sumenq - root.blocks[b].size`.  $\square$

**Corollary 45.** If operations in the root are applied with the order of  $L$ , the number of non-null dequeues in `root.blocks[b]` is `root.blocks[b].numenq - root.blocks[b].size + root.blocks[b - 1].size`.

**Lemma 46.**  $\text{Resp}(D_i(\text{root}, b))$  is null iff `root.blocks[b - 1].size + root.blocks[b].numenq - i` < 0.

*Proof.* From Corollary 45 and Lemma 18.  $\square$

**Lemma 47.** `FindResponse(b, i)` returns  $\text{Resp}(D_i(\text{root}, b))$ .

*Proof.*  $D_i(\text{root}, b)$  is  $D_{\text{root.blocks}[b-1].\text{sum}_{\text{deq}}+i}(\text{root})$  by Definition 15 and Lemma 19.  $D_i(\text{root}, b)$  returns null at Line 220 if `root.blocks[b - 1].size + root.blocks[b].numenq - i` < 0 and  $\text{Resp}(D_i(\text{root}, b)) = \text{null}$  in this case by Lemma 46. Otherwise, if  $D_i(\text{root}, b)$  is the  $i$ th non-null dequeue in  $L$  it should return the  $i$ th enqueued value. By Lemma 44 there are `root.blocks[b - 1].sumenq - root.blocks[b - 1].size` non-null



Dequeue operations in `root.blocks[0..b-1]`. The Dequeues in `root.blocks[b]` before  $D_i(\text{root}, b)$  are non-null dequeues. So  $D_i(\text{root}, b)$  is the  $e$ th non-null Dequeue where  $e = i + \text{root.blocks}[b-1].\text{sum\_deq} - \text{root.blocks}[b-1].\text{size}$  (Line 222). See Figure 20.

After computing  $e$  at Line 222, the code finds  $b, i$  such that  $E_i(\text{root}, b) = E_e(\text{root})$  using `DSearch` and then finds its `element` using `GetEnqueue` (Line 223). Correctness of `DSearch` and `GetEnqueue` routines are shown in Lemmas 30 and 31. □

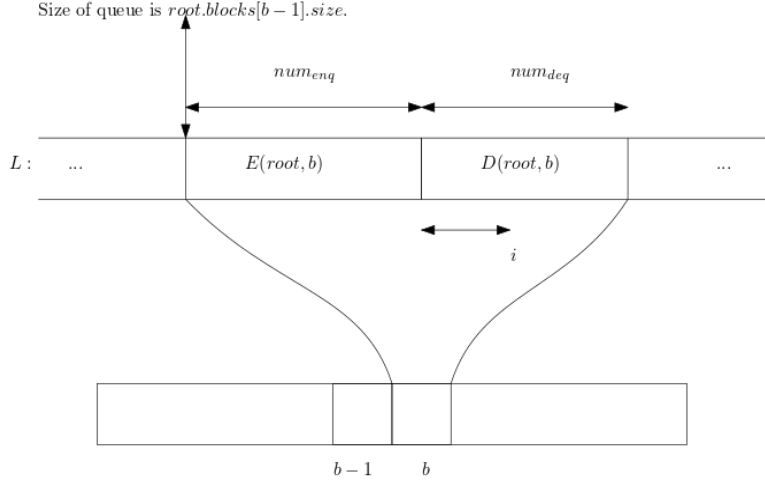


Figure 20: The position of  $D_i(\text{root}, b)$ .

**Lemma 48.** *The responses to operations in our algorithm would be the same as in the sequential execution in the order given by  $L$ .*

*Proof.* **Enqueue** operations do not return any value. By Lemma 47 response of a **Dequeue** in our algorithm is same as the response from the sequential execution of  $L$ . □

**Theorem 49 (Main).** *The queue implementation is linearizable.*

*Proof.* The theorem follows from Lemmas 40 and 48. □

**Remark** In fact our algorithm is strongly linearizable as defined in [5]. By Definition 15 the linearization ordering of operations will not change as blocks containing new operations are appended to the root.

## 5 Analysis

**Lemma 50** (Amortized time analysis). *Enqueue() and Dequeue(), each take  $O(\log^2 p + \log q)$  steps in amortized analysis. Where  $p$  is the number of processes and  $q$  is the size of the queue at the time of invocation of operation.*

*Proof.* **Enqueue(x)** consists of creating a **block(x)** and appending it to the tree. The first part takes constant time. To propagate **x** to the root the algorithm tries two **Refreshes** in each node of the path from the leaf to the root (Lines 302, 303). We can see from the code that each **Refresh** takes constant number of steps since creating a block is done in constant time and does  $O(1)$  CASes. Since the height of the tree is  $\Theta(\log p)$ , **Enqueue(x)** takes  $O(\log p)$  steps.

A **Dequeue()** creates a block with null value element, appends it to the tree, computes its order among enqueue operations, and returns the response. The first two part is similar to an **Enqueue** operation. To compute the order of a **dqueue** in  $D(n)$  there are some constant steps and **IndexDequeue()** is called. **IndexDequeue** does a search with range  $p$  in each level which takes  $O(\log^2 p)$  in the tree. In the **FindResponse()** routine **DSearch()** in the root takes  $\Theta(\log(\text{root.blocks}[b].\text{size} + \text{root.blocks}[\text{end}].\text{size}))$  by Lemma 30, which is  $O(\log \text{ size of the queue when enqueue is invoked} + \log \text{ size of the queue when dequeue is invoked})$ . Each search in **GetEnqueue()** takes  $O(\log p)$  since there are  $\leq p$  subblocks in a block (Lemma 29), so **GetEnqueue()** takes  $O(\log^2 p)$  steps.

If we split **DSearch** time cost between the corresponding **Enqueue**, **Dequeue**, in amortized we have **Enqueue** takes  $O(\log p + q)$  and **Dequeue** takes  $O(\log^2 p + q)$  steps.  $\square$

**Lemma 51.** *An Enqueue() or Dequeue() operation, does at most  $4\log p$  CAS operations.*

*Proof.* In each height of the tree at most 2 times **Refresh** is invoked and every **Refresh** invokes at most 3 CASes, one in Line 320 and two from **Advance** in Line 327.  $\square$

**Lemma 52** (**DSearch** Analysis). *If the element enqueued by  $E_i(\text{root}, b) = E_e(\text{root})$  is the response to some Dequeue() operation in  $\text{root.blocks}[\text{end}]$ , then  $\text{DSearch}(e, \text{end})$  takes  $O(\log(\text{root.blocks}[b].\text{size} + \text{root.blocks}[\text{end}].\text{size}))$  steps.*

*Proof.* First we show  $\text{end} - b - 1 \leq 2 \times \text{root.blocks}[b-1].\text{size} + \text{root.blocks}[\text{end}].\text{size}$ . Suppose there were more than  $\text{root.blocks}[b].\text{size}$  **Dequeues** in  $\text{root.blocks}[b+1 \dots \text{end}-1]$ . Then the element in the queue which is the response to the **Dequeue()** would become dequeued at some point before

`root.blocks[end]`'s first `Dequeue()`. Furthermore in the execution of queue operations in the linearization ordering, the size of the queue becomes `root.blocks[end].size` after the operations of `root.blocks[end]`. There can be at most `root.blocks[b].size` `Dequeues` in `root.blocks[b + 1 .. end - 1]`; otherwise all elements enqueued by `root.blocks[b]` would be dequeued before `root.blocks[end]`. The final size of the queue after `root.blocks[1 .. end]` is `root.blocks[end].size`. After an execution on a queue the *size* of the queue is greater than or equal to  $\#enqueues - \#dequeues$  in the execution. We know the number of dequeues in `root.blocks[b + 1 .. end - 1]` is less than `root.blocks[b].size`, therefore there cannot be more than `root.blocks[b].size + root.blocks[end].size` `Enqueues` in `root.blocks[b + 1 .. end - 1]`. Overall there can be at most  $2 \times \text{root.blocks}[b].\text{size} + \text{root.blocks}[end].\text{size}$  operations in `root.blocks[b + 1 .. end]` and since from line 318 we know that `num` field of the every block in the tree is greater than 0, each block has at least one operation, there are at most  $2 \times \text{root.blocks}[b].\text{size} + \text{root.blocks}[end].\text{size}$  blocks in between `root.blocks[b]` and `root.blocks[end]`. So  $end - b - 1 \leq 2 \times \text{root.blocks}[b].\text{size} + \text{root.blocks}[end].\text{size}$ .

So the doubling search reaches `start` such that the `root.blocks[start].sumenq` is less than  $e$  in  $O(\log(\text{root.blocks}[b].\text{size} + \text{root.blocks}[end].\text{size}))$  steps. See Figure 21. After Line 805, the binary search that finds  $b$  also takes  $O(\log(\text{root.blocks}[b].\text{size} + \text{root.blocks}[end].\text{size}))$ . Next,  $i$  is computed via the definition of `sumenq` in constant time (Line 807). So the claim is proved.  $\square$

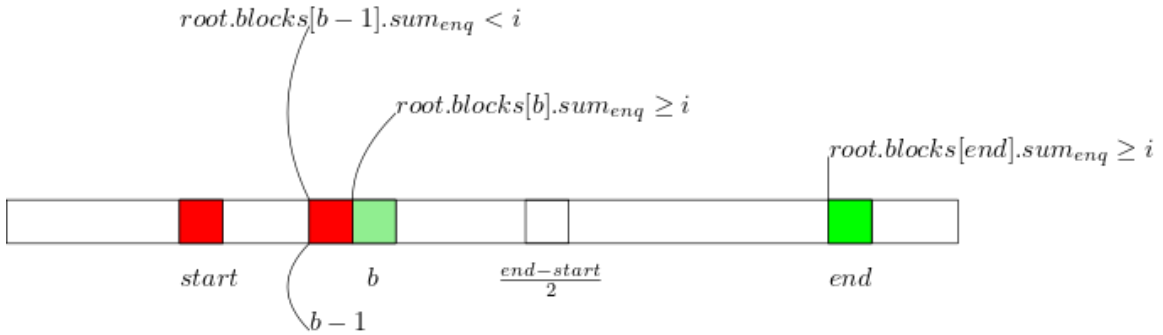


Figure 21: Distance relations between `start`,  $b$ , `end`.

## 5.1 Garbage Collection or Getting rid of the infinite Arrays

## 6 Using Queues to Implement Vectors

Supporting Append, Read, Write in PolyLog time by modifying Get(Enq) Method. Create a Universal Construction Using our vector

## 7 Conclusion

possible directions for work

Maybe Stacks

Characterize what datastructure can be used for this approach, we already know: queue, fetch & Inc, Vectors

## References

## References

- [1] Tushar Deepak Chandra, Prasad Jayanti, and King Tan. A polylog time wait-free construction for closed objects. In Brian A. Coan and Yehuda Afek, editors, *Proceedings of the Seventeenth Annual ACM Symposium on Principles of Distributed Computing, PODC '98, Puerto Vallarta, Mexico, June 28 - July 2, 1998*, pages 287–296. ACM, 1998.
- [2] Robert Colvin and Lindsay Groves. Formal verification of an array-based nonblocking queue. In *10th International Conference on Engineering of Complex Computer Systems (ICECCS 2005), 16-20 June 2005, Shanghai, China*, pages 507–516. IEEE Computer Society, 2005.
- [3] Faith Ellen and Philipp Woelfel. An optimal implementation of fetch-and-increment. In *Proceedings of the 27th International Symposium on Distributed Computing - Volume 8205, DISC 2013*, page 284298, Berlin, Heidelberg, 2013. Springer-Verlag.
- [4] Anders Gidenstam, Håkan Sundell, and Philippas Tsigas. Cache-aware lock-free queues for multiple producers/consumers and weak memory consistency. In Chenyang Lu, Toshimitsu Masuzawa, and Mohamed Mosbah, editors, *Principles of Distributed Systems - 14th International Conference, OPODIS 2010, Tozeur, Tunisia, December 14-17, 2010. Proceedings*, volume 6490 of *Lecture Notes in Computer Science*, pages 302–317. Springer, 2010.
- [5] Wojciech M. Golab, Lisa Higham, and Philipp Woelfel. Linearizable implementations do not suffice for randomized distributed computation. In Lance Fortnow and Salil P. Vadhan, editors, *Proceedings of the 43rd ACM Symposium on Theory of Computing, STOC 2011, San Jose, CA, USA, 6-8 June 2011*, pages 373–382. ACM, 2011.

- [6] Danny Hendler, Itai Incze, Nir Shavit, and Moran Tzafrir. Flat combining and the synchronization-parallelism tradeoff. In Friedhelm Meyer auf der Heide and Cynthia A. Phillips, editors, *SPAA 2010: Proceedings of the 22nd Annual ACM Symposium on Parallelism in Algorithms and Architectures, Thira, Santorini, Greece, June 13-15, 2010*, pages 355–364. ACM, 2010.
- [7] Maurice Herlihy. Wait-free synchronization. *ACM Trans. Program. Lang. Syst.*, 13(1):124149, jan 1991.
- [8] Moshe Hoffman, Ori Shalev, and Nir Shavit. The baskets queue. In Eduardo Tovar, Philippas Tsigas, and Hacène Fouchal, editors, *Principles of Distributed Systems, 11th International Conference, OPODIS 2007, Guadeloupe, French West Indies, December 17-20, 2007. Proceedings*, volume 4878 of *Lecture Notes in Computer Science*, pages 401–414. Springer, 2007.
- [9] Prasad Jayanti. A time complexity lower bound for randomized implementations of some shared objects. In Brian A. Coan and Yehuda Afek, editors, *Proceedings of the Seventeenth Annual ACM Symposium on Principles of Distributed Computing, PODC '98, Puerto Vallarta, Mexico, June 28 - July 2, 1998*, pages 201–210. ACM, 1998.
- [10] Prasad Jayanti and Srdjan Petrovic. Logarithmic-time single deleter, multiple inserter wait-free queues and stacks. In Ramaswamy Ramanujam and Sandeep Sen, editors, *FSTTCS 2005: Foundations of Software Technology and Theoretical Computer Science, 25th International Conference, Hyderabad, India, December 15-18, 2005, Proceedings*, volume 3821 of *Lecture Notes in Computer Science*, pages 408–419. Springer, 2005.
- [11] Alex Kogan and Erez Petrank. Wait-free queues with multiple enqueueers and dequeuers. In Calin Cascaval and Pen-Chung Yew, editors, *Proceedings of the 16th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, PPOPP 2011, San Antonio, TX, USA, February 12-16, 2011*, pages 223–234. ACM, 2011.
- [12] Edya Ladan-Mozes and Nir Shavit. An optimistic approach to lock-free FIFO queues. *Distributed Comput.*, 20(5):323–341, 2008.
- [13] Maged M. Michael and Michael L. Scott. Simple, fast, and practical non-blocking and blocking concurrent queue algorithms. In James E. Burns and Yoram Moses, editors, *Proceedings of the Fifteenth Annual ACM Symposium on Principles of Distributed Computing, Philadelphia, Pennsylvania, USA, May 23-26, 1996*, pages 267–275. ACM, 1996.

- [14] Mark Moir, Daniel Nussbaum, Ori Shalev, and Nir Shavit. Using elimination to implement scalable and lock-free FIFO queues. In Phillip B. Gibbons and Paul G. Spirakis, editors, *SPAA 2005: Proceedings of the 17th Annual ACM Symposium on Parallelism in Algorithms and Architectures, July 18-20, 2005, Las Vegas, Nevada, USA*, pages 253–262. ACM, 2005.
- [15] Adam Morrison and Yehuda Afek. Fast concurrent queues for x86 processors. In Alex Nicolau, Xiaowei Shen, Saman P. Amarasinghe, and Richard W. Vuduc, editors, *ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, PPOPP '13, Shenzhen, China, February 23-27, 2013*, pages 103–112. ACM, 2013.
- [16] Niloufar Shafiei. Non-blocking array-based algorithms for stacks and queues. In Vijay K. Garg, Roger Wattenhofer, and Kishore Kothapalli, editors, *Distributed Computing and Networking, 10th International Conference, ICDCN 2009, Hyderabad, India, January 3-6, 2009. Proceedings*, volume 5408 of *Lecture Notes in Computer Science*, pages 55–66. Springer, 2009.
- [17] Philippas Tsigas and Yi Zhang. A simple, fast and scalable non-blocking concurrent FIFO queue for shared memory multiprocessor systems. In Arnold L. Rosenberg, editor, *Proceedings of the Thirteenth Annual ACM Symposium on Parallel Algorithms and Architectures, SPAA 2001, Heraklion, Crete Island, Greece, July 4-6, 2001*, pages 134–143. ACM, 2001.