

Wait-free Queues with Polylogarithmic Step Complexity

Hossein Naderibeni

supervised by Eric Ruppert

September 27, 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

1	Introduction	2
2	Related Work	4
2.1	List-based Queues	4
2.2	Universal Constructions	5
2.3	Attiya Fourier Lower Bound	5
3	Our Queue	6
3.1	Pseudocode description	11
3.2	Pseudocode	13
4	Proof of Correctness	17
4.1	Garbage Collection or Getting rid of the infinite Arrays	30
5	Using Queues to Implement Vectors	31
6	Conclusion	32

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.

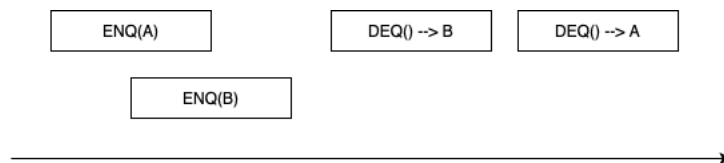


Figure 1: An example of a linearizable execution. Either `Enqueue(A)` or `Enqueue(B)` could take effect first since they are concurrent.

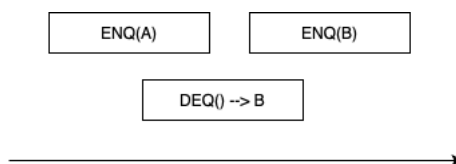


Figure 2: An example of an execution that is not linearizable. Since `Enqueue(A)` has completed before `Enqueue(B)` is invoked the `Dequeue()` should return A or nothing.

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 [?] proved an $\Omega(\log p)$ lower bound on the worst-case shared-access time complexity of p -process universal constructions. He also introduced [?] a construction that achieves $O(\log^2 p)$ shared accesses. Here, we first introduce a universal construction using $O(\log p)$ CAS operations [?]. 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 [?] 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.

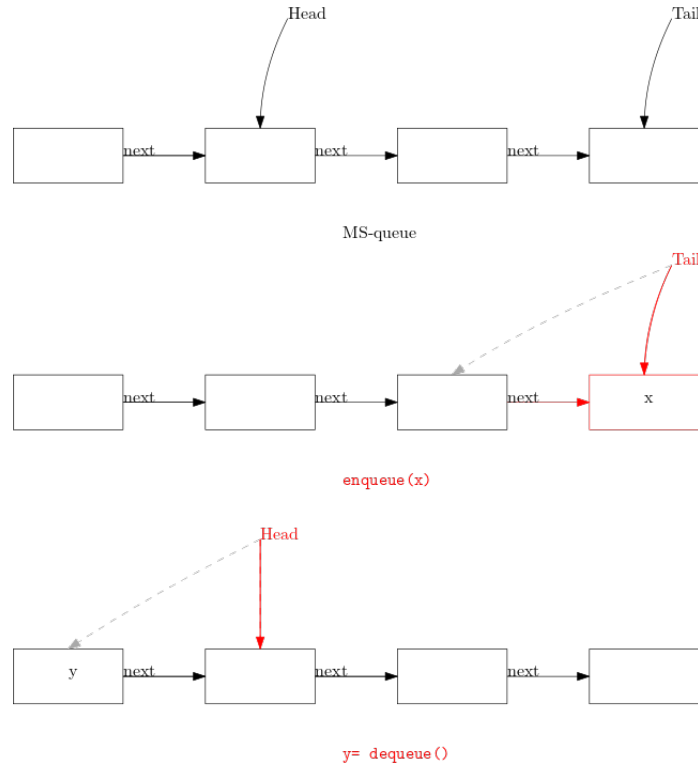


Figure 3: MS-queue structure, enqueue and dequeue operations. In the first diagram the first element has been dequeued. Red arrows show new pointers and gray dashed arrows show the old pointers.

Moir, Nussbaum, and Shalev [?] 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 [?] 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.

Ladan-Mozes and Shavit [?] 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. [?] proposed a new paradigm called flat combining. Their queue is linearizable but not lock-free. Their main idea is

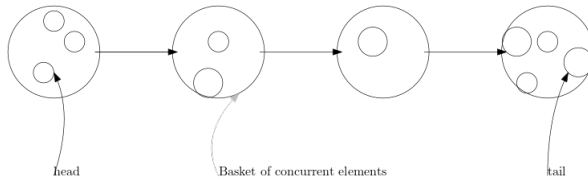


Figure 4: Baskets queue idea. There is a time that all operations in a basket were running concurrently, but only one has succeeded to do CAS. To order the operations in a basket, the mechanism in the algorithm for processes is to CAS again. The successful process will be the next one in the basket and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

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 [?] 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 $\Omega(p)$.

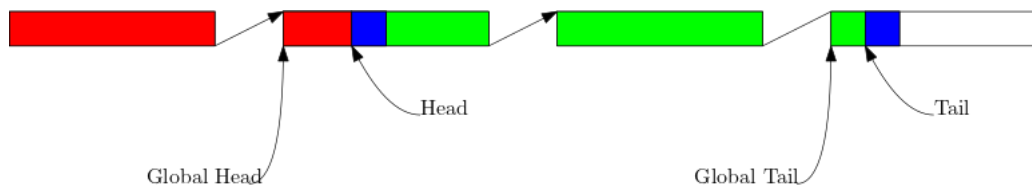


Figure 5: Global pointers point to arrays. Head and Tail elements are blue, dequeued elements are red and current elements of the queue are green.

Kogan and Petrank [?] 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* [?]. It is not limited to list-based queues and array-based queues share the CAS retry problem as well [?, ?, ?]. 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 [?]. 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 [?]. He also introduced a construction that achieves $O(\log^2 p)$ shared accesses [?]. 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 [?]. 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-dequeuer queue [?]. We benefit from some ideas of their work to design a polylogarithmic multi-enqueuer multi-dequeuer 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.

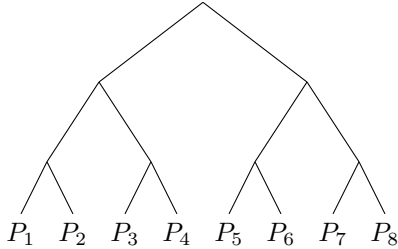


Figure 6: Each process has a leaf and in each node there is an ordering of operations stored. Each node tries to propagate its operations up to the root, which stores the total ordering of all operations.

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.

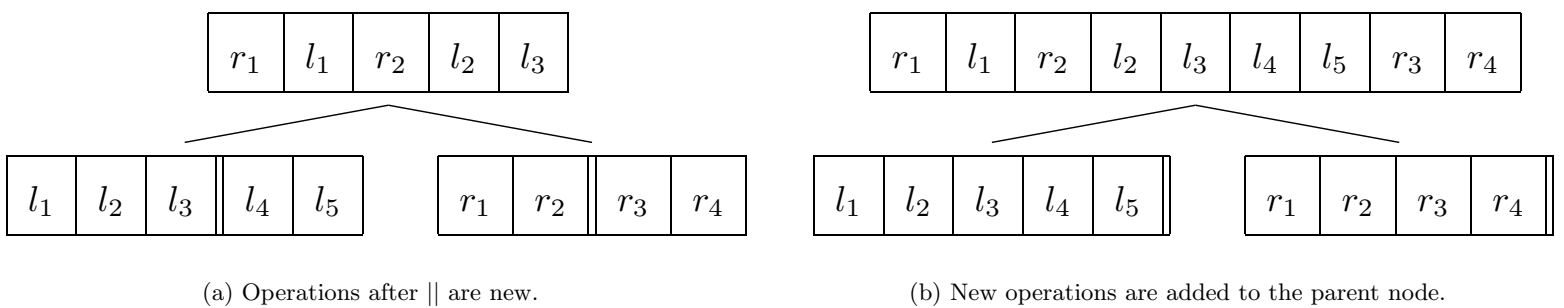


Figure 7: Before and after of a `REFRESH(n)` with successful **CAS**. Operations propagating from the left child are numbered with l and from the right child by r and the operations in children after `||` are new.

$fx ||$

The solution for implementing the orderings in the tree told above is not efficient, because there are big **CAS**es 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 `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

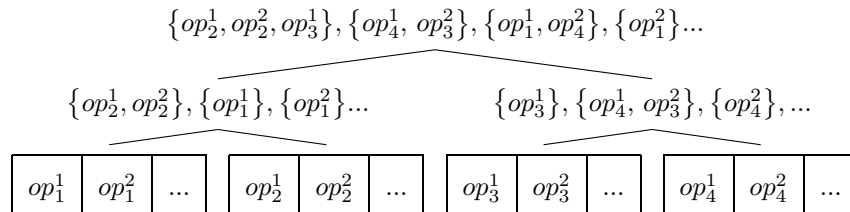


Figure 8: In each internal node, we store the set of all the operations propagated together, and one can arbitrarily linearize the sets of concurrent operations among themselves. Since we linearize operations when they are added to the root, ordering the blocks in the root is important.

invoke two operations concurrently. Furthermore since the operations in a `REFRESH()` step are concurrent we can linearize them among themselves in any order we wish. Note that if two operations are in read one `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 `GetENQ(n,b,i)`. `DSearch(i)` finds the block containing the i th enqueue in the root and `GetENQ(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 ?? we know a block size is at most p . So `BinarySearch` takes at most $(O)(\log p)$, since with knowing pointers of a block and its previous block we can determine the base (domain ?) to search and its size is $O(p)$.

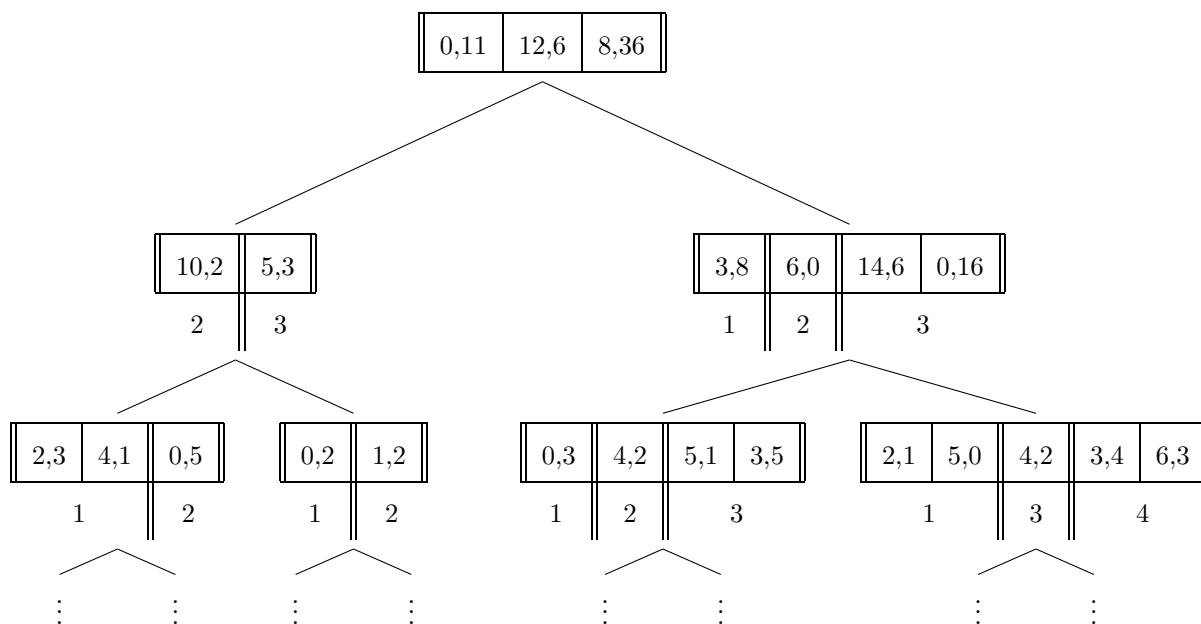


Figure 9: Showing concurrent operation sets with blocks. Each block consists of a pair(left, right) indicating the number of operations from the left and the right child, respectively. Block (12,6) in the root contains blocks (10,2) from the left child and (6,0) from the right child. Blocks between two 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.

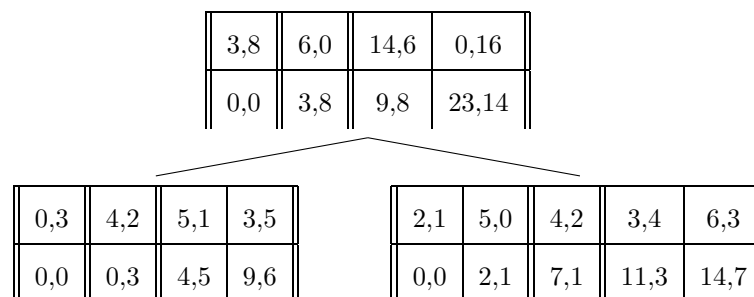


Figure 10: Using Prefix sums in blocks. When we want to find block b elements in its children, we can use binary search. The number below each block shows the count of elements in the previous blocks.

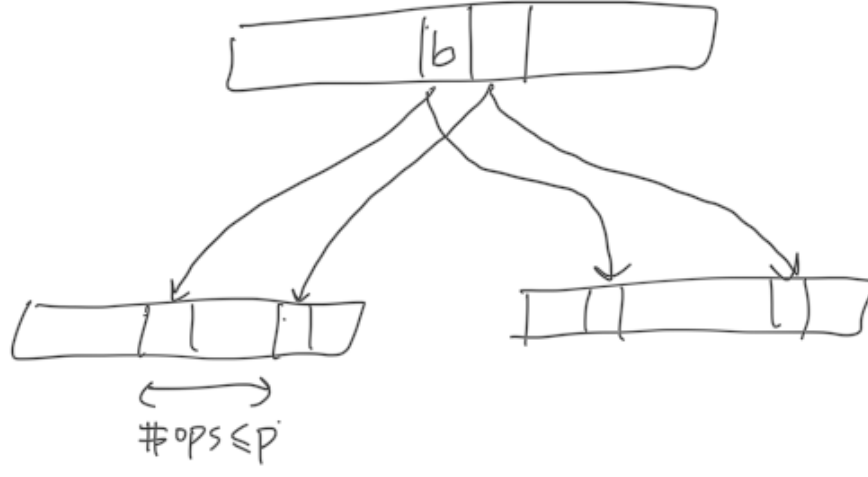


Figure 11: Block have pointers to the starting block of theirs for each child.

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.

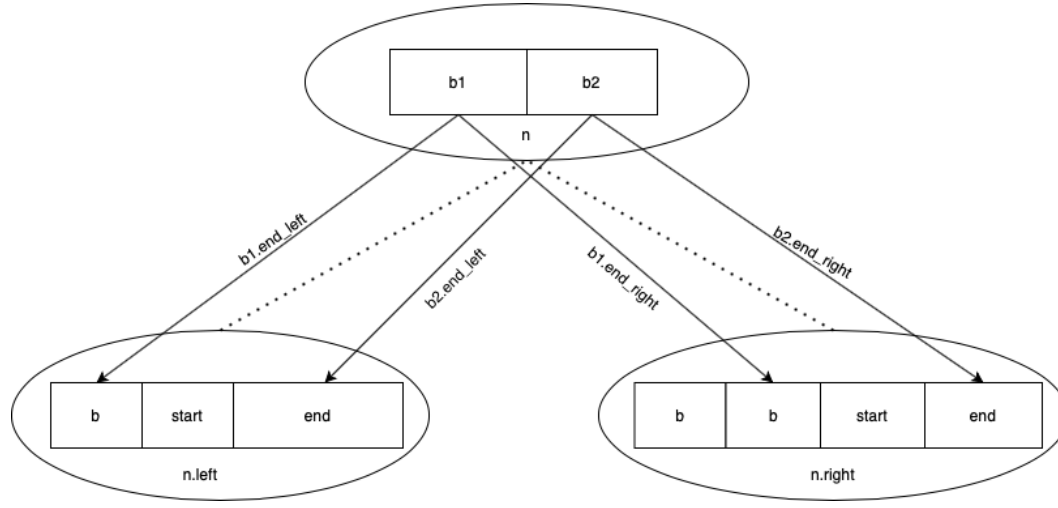


Figure 12: Snapshot of a CreateBlock()

Computing $Get(n, b, i)$

How Refresh(n) works.

Computing superblock

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

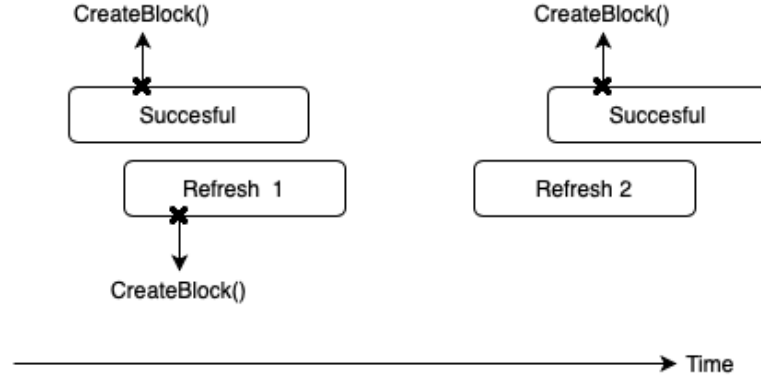


Figure 13: The second failed Refresh is assuredly concurrent to a Successful Refresh() with CreateBlock line after first failed Refresh's CreateBlock().

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 sufficient 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 ?? 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.

DEQ()	ENQ(5), ENQ(2), ENQ(1), DEQ()	ENQ(3), DEQ()	ENQ(4), DEQ(), DEQ(), DEQ(), DEQ()
-------	-------------------------------	---------------	------------------------------------

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.

	DEQ()	ENQ(5), ENQ(2), ENQ(1), DEQ()	ENQ(3), DEQ()	ENQ(4), DEQ(), DEQ(), DEQ(), DEQ()
#enqueues	0	3	1	1
#dequeues	1	1	1	4
#non-null dequeues	0	1	2	5
size	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:

$\text{size}[b] = \max(\text{size}[b-1] + \text{enqueues}[b] - \text{dequeues}[b], 0)$

$\text{non-null dequeues}[b] = \text{non-null dequeues}[b-1] + \text{dequeues}[b] - \text{size}[b-1] - \text{enqueues}[b]$

Given DEQ is in block b , $\text{response}(\text{DEQ})$ would be:

$(\text{size}[b-1] - \text{index of DEQ in the block's dequeues} \geq 0) ? \text{ENQ}[\text{non-null dequeues}[b-1] + \text{index of DEQ in the block's dequeues}]$
: null;

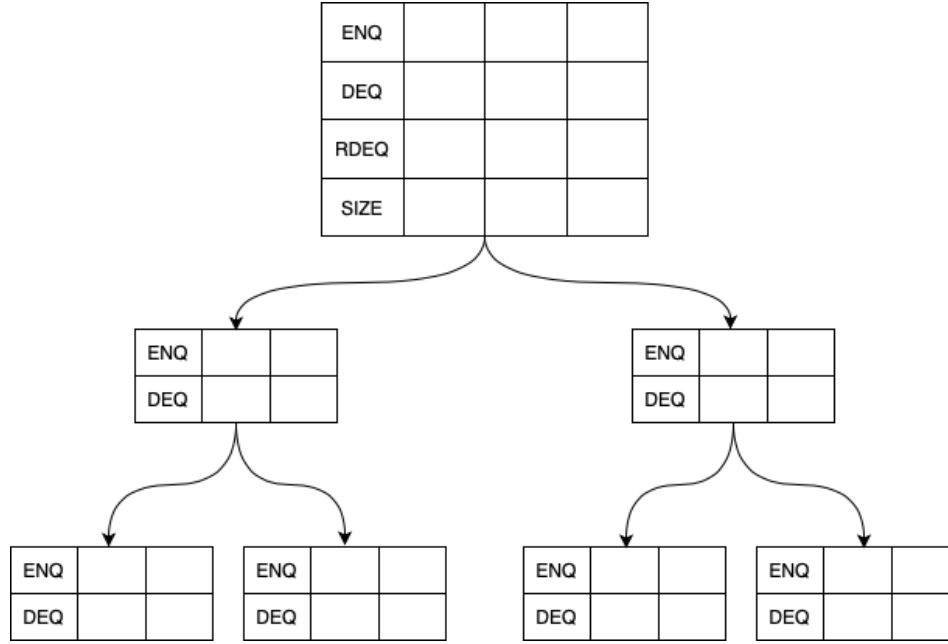


Figure 14: Fields stored in the Queue nodes.

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 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 `numenq-left`, `numdeq-left`, `numenq-right`, `numdeq-right`. 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 `sumnon-null deq`. `size` in a **block** shows the number of elements after the **block** has finished and `sumnon-null deq` 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

► 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_{enq}, sum_{deq} : # enqueue, dequeue operations in the prefix for the block

► InternalBlock extends Block

- *int* end_{left}, end_{right} : indices of the last subblock of the block in the left and right child
- *int* sum_{enq-left} : # enqueue operations in the prefix for left.blocks[end_{left}]
- *int* sum_{deq-left} : # dequeue operations in the prefix for left.blocks[end_{left}]
- *int* sum_{enq-right} : # enqueue operations in the prefix for right.blocks[end_{right}]
- *int* sum_{deq-right} : # dequeue operations in the prefix for right.blocks[end_{right}]

► RootBlock extends InternalBlock

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

Abbreviations:

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

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, computes its order among operations, 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>= INDEXDEQ(leaf.head, 1)
215:   output= FINDRESPONSE(b, i)
216:   return output
217: end DEQUEUE

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

Algorithm Root

```
Precondition: root.blocks[end].sumenq ≥ e

801: <int, int> DSEARCH(int e, int end) ▷ Returns <b,i> if  $E_e(root) = E_i(root, b)$ .
802:   start= end-1
803:   while root.blocks[start].sumenq ≥ e do
804:     start= max(start-(end-start), 0)
805:   end while
806:   b= root.BinarySearch(sumenq, e, start, end)
807:   i= e- root.blocks[b-1].sumenq
808:   return <b,i>
809: end DSEARCH
```

Algorithm Leaf

```
601: void APPEND(block blk) ▷ Append is only called by the owner of the leaf.
602:   blocks[head]= blk
603:   head+=1
604:   parent.PROPAGATE()
605: end APPEND
```

Algorithm Node

```
301: void PROPAGATE() ↔ Precondition: blocks[start..end] contains a block with field f ≥ i
302:   if not REFRESH() then
303:     REFRESH() ▷ Does binary search for the value
304:   end if i of the given prefix sum field. Returns the index of the leftmost block in
305:   if this is not root then blocks[start..end] whose field f is ≥ i.
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, h)
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:   hp= parent.head
322:   this.ADVANCE(h, hp)
323:   return result
324: end REFRESH

325: void ADVANCE(int h, int hp)
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: end BINARYSEARCH

331: <Block, int, int> CREATEBLOCK(int i) ▷ Creates and returns the block
    to be inserted as ith block in blocks.
332:   block newBlock= NEW(block)
333:   for each dir in {left, right} do
334:     indexlast= dir.head-1
335:     indexprev= blocks[i-1].enddir
336:     newBlock.enddir= indexlast
337:     blocklast= dir.blocks[indexlast]
338:     blockprev= dir.blocks[indexprev]
339:     ▷ newBlock includes dir.blocks[indexprev+1..indexlast].
340:     newBlock.sumenq-dir= blocks[i-1].sumenq-dir + blocklast.sumenq
    - blockprev.sumenq
341:     newBlock.sumdeq-dir= blocks[i-1].sumdeq-dir + blocklast.sumdeq
    - blockprev.sumdeq
342:   end for
343:   if this is root then
344:     newBlock.size = max(root.blocks[i-1].size + newBlock.numenq
    - newBlock.numdeq, 0)
345:   end if
346:   return <b, nleft, nright>
347: end CREATEBLOCK
```

Algorithm Node

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

401: *element* GETENQ(*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(sum_{enq}, i+blocks[b-1].sum_{enq-left}, blocks[b-1].end_{left}+1, blocks[b].end_{left})

406: **return** left.GETENQ(subBlock, i)

407: **else**

408: i= i-blocks[b].num_{enq-left}

409: subBlock= right.BINARYSEARCH(sum_{enq}, i+right.blocks[b-1].sum_{enq-right}, blocks[b-1].end_{right}+1, blocks[b].end_{right})

410: **return** right.GETENQ(subBlock, i)

411: **end if**

412: **end** GETENQ

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

413: <*int*, *int*> INDEXDEQ(*int* b, *int* i) ▷ Returns <x, y> if $D_{\text{this}, b, i} = D_{\text{root}, x, y}$.

414: **if** this is root **then**

415: **return** <b, i>

416: **else**

417: dir= (parent.left==n)? left: right ▷ check if this node is a left or a right child

418: superBlock= parent.BINARYSEARCH(sum_{deq-dir}, i+blocks[b-1].sum_{deq}, blocks[b].super, blocks[b].super+1)

▷ superblock's group has at most p difference with the value stored in **super** [].

419: **if** dir is left **then**

420: i+= blocks[b-1].sum_{enq}-blocks[superBlock-1].sum_{enq-left} ▷ consider the enqueues in the previous blocks from the left child

421: **end if**

422: **if** dir is right **then**

423: i+= blocks[b-1].sum_{enq}-blocks[superBlock-1].sum_{enq-right} ▷ consider the enqueues in the previous blocks from the right child

424: i+= blocks[superBlock].num_{deq-left} ▷ consider the dequeues from the right child

425: **end if**

426: **return** this.parent.INDEXDEQ(superBlock, i)

427: **end if**

428: **end** INDEXDEQ

4 Proof of Correctness

To prove a shared data structure works correctly, it is sufficient to show it is linearizable. 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 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 operation of each block which we use to prove the correctness of `DSearch()`, `GetEnq()` and `IndexDeq()`. Finally (2) is followed by proving the correctness of the way we compute the response of a dequeue.

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 $n.blocks[i]$ and b' be $n.blocks[j]$. 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 b to be the blocks in $n.blocks[0..i]$.

Next, we show that the value of `head` in a node can be increased. After 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 , $n.head$ 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 node n . After R terminates, $n.head$ 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 $n.head$ has changed from the value that was read in Line 310. By Observation 2 this means another process has incremented $n.head$. □

Now we show $n.blocks[n.head]$ is the last block written into n or the first empty block in n .

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

Proof. Initially the invariant is true since $n.head = 1$, $n.blocks[0] \neq \text{null}$ and $n.blocks[x] \neq \text{null}$ for every $x > 0$. The truth of the invariant may be affected by writing into $n.blocks$ or incrementing $n.head$. We show the invariant still holds after these two changes.

In the algorithm, $n.blocks$ is modified only on Line 320, which updates $n.blocks[h]$ where h is the value read from $n.head$ in Line 310. Since the CAS in Line 320 is successful it means $n.head$ has not changed from h before doing the CAS, because if so then $n.blocks[h].CAS$ could not be successful. Writing into $n.blocks[n.head = h]$ preserves the invariant, since the claim does not talk about the content of $n.blocks[n.head]$.

The value of $n.head$ is modified only in Line 327. If $n.head$ is incremented to $h + 1$ it is sufficient to show $n.blocks[h] \neq \text{null}$. `Advance()` is called in Lines 322 and 314. In case `advance`, Line 320 was finished before doing 327 whether 320 is successful or not $n.blocks[h] \neq \text{null}$. In the other case also $n.blocks[h] \neq \text{null}$ because of the `if` condition in Line 313. □

We define the subblocks of a block recursively.

Definition 5 (Subblock). A block is a *direct subblock* of i th block in node n if it is in $n.left.blocks[n.blocks[i-1].end_{left}+1..n.blocks[i].end_{left}]$ or in $n.right.blocks[n.blocks[i-1].end_{right}+1..n.blocks[i].end_{right}]$. 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 $n.blocks$ or is a subblock of a block in $n.blocks$.

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 336 of `CreateBlock(i)` was one less than the value read at Line 334 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 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 then 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 $[n.\text{blocks}[i-1].\text{end}_{\text{left}} + 1 \cdots n.\text{blocks}[i].\text{end}_{\text{left}}]$ and $[n.\text{blocks}[j-1].\text{end}_{\text{left}} + 1 \cdots n.\text{blocks}[j].\text{end}_{\text{left}}]$ cannot overlap. Therefore, direct subblocks of $n.\text{blocks}[i]$ and $n.\text{blocks}[j]$ cannot overlap. \square

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

Definition 8 (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 set of the operations in leaf subblocks of b . We denote the set of operations of block b by $\text{ops}(b)$ and the union set of operations of set of blocks B by $\text{ops}(B)$. We also say b contains op if $op \in \text{ops}(b)$.

Definition 9 (Superblock). Block b is *superblock* of block c if c is a direct subblock of b .

Corollary 10. *Every block has at most one superblock.*

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

Operations are distinct `Enqueues` and `Dequeues` invoked by processes. 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 $n.\text{blocks}[i]$ and $n.\text{blocks}[j]$. 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 leads us to contradictory with the hypothesis. \square

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

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

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

Proof. Blocks are only appended (not modified) with CAS to $n.\text{blocks}[n.\text{head}]$, 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

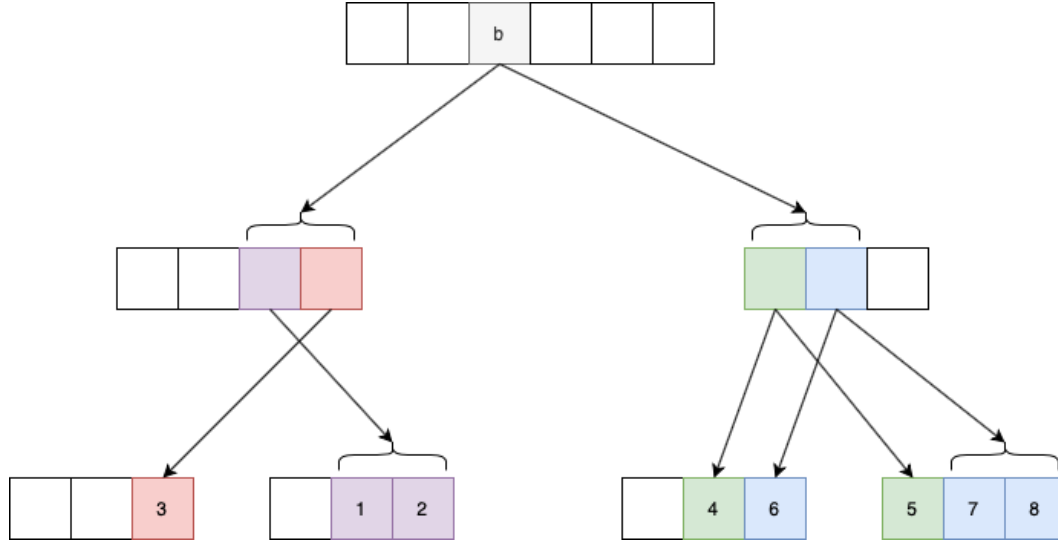


Figure 15: Order of operations in b . Operations in the leaves are ordered with numerical order shown in the drawing.

Now we define the ordering of operations stored in each node. In the non-root nodes we only need ordering of 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 25 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 $leaf.blocks[b].num_{enq} = 0$. If n is an internal node, then

$$E(n, b) = E(n.left, n.blocks[b-1].end_{left} + 1) \cdot E(n.left, n.blocks[b-1].end_{left} + 2) \cdots E(n.left, n.blocks[b].end_{left}) \cdot \\ E(n.right, n.blocks[b-1].end_{right} + 1) \cdot E(n.right, n.blocks[b-1].end_{right} + 2) \cdots E(n.right, n.blocks[b].end_{right})$$

- $E_i(n, b)$ 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 $leaf.blocks[b].num_{deq} = 0$. If n is an internal node, then
- $$D(n, b) = D(n.left, n.blocks[b-1].end_{left} + 1) \cdot D(n.left, n.blocks[b-1].end_{left} + 2) \cdots D(n.left, n.blocks[b].end_{left}) \cdot \\ D(n.right, n.blocks[b-1].end_{right} + 1) \cdot D(n.right, n.blocks[b-1].end_{right} + 2) \cdots D(n.right, n.blocks[b].end_{right})$$
- $D_i(n, b)$ is the i th dequeue in $D(n, b)$.
 - The order of the dequeue operations in the node n : $D(n) = D(n, 1) \cdot D(n, 2) \cdot D(n, 3) \cdots$
 - $D_i(n)$ is the i th dequeue in $D(n)$.

Definition 16 (Linearization). $L = E(root, 1).D(root, 1).E(root, 2).D(root, 2).E(root, 3).D(root, 3) \cdots$

Definition 17. Let t^{op} be the time op is invoked and ^{op}t be the time op terminates. Let t_l be the time immediately before executing Line l and $_lt$ be the time immediately after executing Line l . Let $^{op}_lt$ is the time immediately before running Line l of operation op and t_l^{op} is the immediate time after running Line l of operation op . 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 18 (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 established in the parent after the **Refresh**'s succesful **CAS** at Line 320.

Lemma 19. *If R is a successful instance of $n.\text{Refresh}()$, 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 $EST_{n.\text{left}}^{t^R} = ops(n.\text{left.blocks}[0..n.\text{left.head}_{309} - 1]) \subseteq ops(n.\text{blocks}_{320}) = ops(n.\text{blocks}[0..n.\text{head}_{320}])$. Line 320 stores a block **new** in n that has $\text{end}_{\text{left}} = n.\text{left.head}_{334} - 1$. Therefore by Definition 5, after the successful **CAS** in Line 320 we know all blocks in $n.\text{left.blocks}[1..n.\text{left.head}_{334} - 1]$ are subblocks of $n.\text{blocks}[1..n.\text{head}_{310}]$. Because of Lemma 2 we have $n.\text{left.head}_{309} - 1 < n.\text{left.head}_{334} - 1$ 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

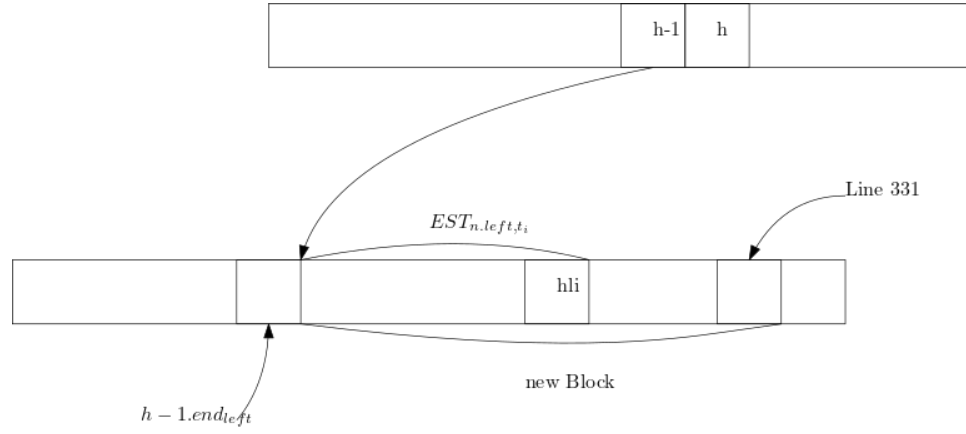


Figure 16: New established operations of the left child are in the new block. (TO UPDATE)

Corollary 20. *If R is a complete instance $n.\text{Refresh}()$, then we have $EST_{n.\text{left}}^{t^R} \cup EST_{n.\text{right}}^{t^R} \subseteq EST_n^{t^R}$.*

Proof. The left hand side is the same as Lemma 19, so it is sufficient to show when R terminates the established blocks in n are a superset of $n.\text{blocks}$ after Line 320. Because of Lemma 3 we are sure that $n.\text{head}$ is incremented after line 327. So the block **new** appended to n at Line 320 is established at Rt . \square

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

Lemma 21. *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 $EST_{n.\text{left}}^{t_{R_1}} \cup EST_{n.\text{right}}^{t_{R_1}} \subseteq EST_n^{t_{R_2}}$.*

Proof. Let R_1 read i from $n.\text{head}$ at Line 310. Note that by Lemma 3 R_1 and R_2 both cannot read the same value i . By Observation 2 R_2 reads 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 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 17). By Lemma 20 we have $EST_{n.\text{left}}^{t_{310}^{R'_2}} \cup EST_{n.\text{right}}^{t_{310}^{R'_2}} \subseteq ops(n.\text{blocks}_{t_{310}^{R'_2}})$. 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 13 $ops(EST_{n.\text{left}, t_{R_1}}) \cup ops(EST_{n.\text{right}, t_{R_1}}) \subseteq ops(n.\text{blocks}_{t_{R_2}})$

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 $t^{R_1} < {}^{R_3}t < {}^{R_3}t < t_{320}^{R_3} < t_{310}^{R_2} < {}^Rt$. From Observation 13 on the operations before and after R_3 's CAS and Lemma 19 on R_3 the claim holds. \square

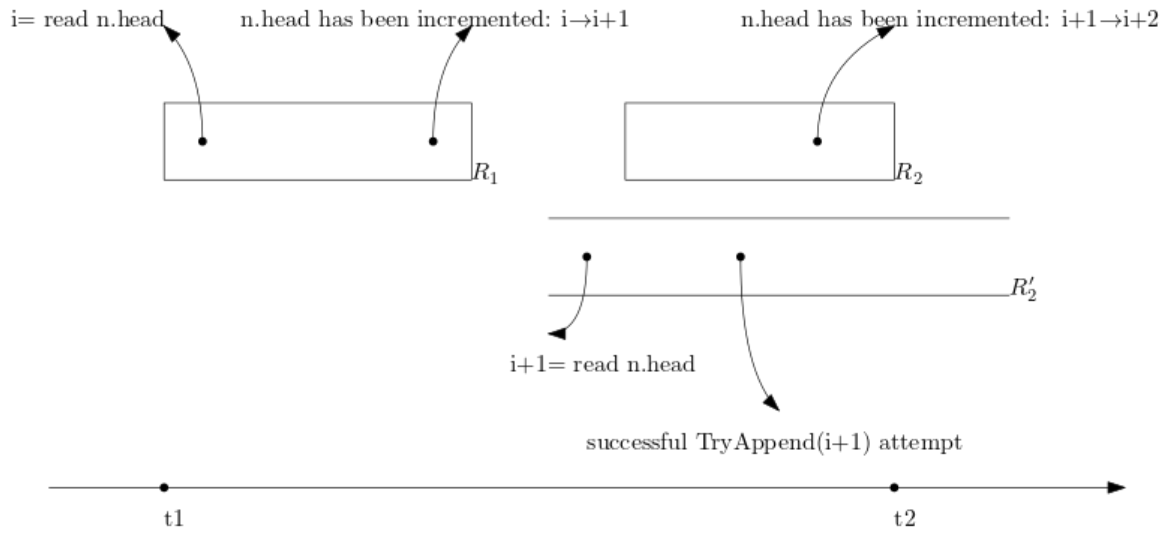


Figure 17: ${}^{R_1}t < t_{310}^{R_1} < \text{incrementing } n.\text{head from } i \text{ to } i+1 < t_{310}^{R'_2} < t_{320}^{R'_2} < \text{incrementing } n.\text{head from } i+1 \text{ to } i+2 < t_{R_2}$

Corollary 22. $ops(EST_{n.\text{left}, t_{302}}) \cup ops(EST_{n.\text{right}, t_{302}}) \subseteq ops(EST_n, t_{303})$

Proof. If the first **Refresh()** in line 302 returns **true** then by Lemma 20 the claim holds. If the first **Refresh()** failed and the second **Refresh()** succeeded the claim still holds by Lemma 20. Otherwise both failed and the claim is satisfied by Lemma 21. \square

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

Corollary 23. *For $A = \text{Append}(b)$ we have $ops(b) \subseteq ops(EST_{n, t^A})$ where $n \in \{\text{nodes in the path from the leaf to the root}\}$.*

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 by Corollary 22 and then calls $n.\text{parent}.\text{Propagate}()$ until n is **root**. \square

Corollary 24. *After `Append(b)` finishes, b is in the nodes the path from the leaf to the root for exactly one time.*

Proof. By the previous corollary and Lemma 24 there is exactly one block in each node containing b . □

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

Lemma 25. *Each block contains at most one operation of each process.*

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 operations concurrently, so op_1 has to be finished before op_2 begins. By Corollary 24 before op_2 calls `Append()`, op_1 exists in every node of the tree on the path from p 's leaf to the root. 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. □

Lemma 26. *Each block has at most p direct subblocks.*

Proof. The claim follows directly from Lemma 25 and the observation that each block appended to an initial 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. □

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

Lemma 27 (DSearch Correctness). *If `root.blocks[end] ≠ null` and $1 \leq e \leq \text{root.blocks[end].sum}_{\text{enq}}$, `DSearch(e, end)` returns $\langle b, i \rangle$ such that $E_i(\text{root}, b) = E_e(\text{root})$.*

Proof. `DSearch` performs a doubling search from `root.blocks[end]` to `root.blocks[0]` to find $E_e(\text{root})$. From Lines 340, 341 we know `sumenq-left`, `sumenq-right` fields of `blocks` in each node are sorted in non-decreasing order. Since `sumenq = sumenq-left + sumenq-right`, `sumenq` values of `root.block[0..end]` is 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`. Block `root.blocks[b]` contains $E_i(\text{root}, b)$. The doubling search on Lines 802–805 doubles its search range in Line 804 and will eventually reach `start` such that `root.blocks[start].sumenq ≤ e ≤ root.blocks[end].sumenq`. In Line 806 Binary Search finds b in the range mentioned. Finally i , is computed from the definition of `sumenq`. □

Lemma 28 (Get correctness). *If $n.\text{blocks}[b].\text{num}_{\text{enq}} \geq i$ then $n.\text{GetENQ}(b, i)$ 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, n is a leaf. Leaf blocks each contain exactly one operation, so only $n.\text{GetENQ}(b, 1)$ can be called when n is a leaf. Line 403 of $n.\text{GetENQ}(b, 1)$ returns the `element` of the `enqueue` operation stored in the b th block of leaf n .

For the induction step we prove if $n.\text{child}.\text{GetENQ}(sb, i)$ returns $E_i(n.\text{child}, sb)$ then $n.\text{GetENQ}(b, i)$ returns $E_i(n, b)$. From Definition 15 of $E(n, b)$ operations from the left subblocks come before the operations from the right subblocks in a block (See Figure 18). `numenq-left` field in $n.\text{blocks}[b]$ stores the number of `Enqueue()` operations from the blocks's subblocks in the left child of n . So the i th enqueue operation is propagated from the right child if and only if i is greater than $b.\text{num}_{\text{enq-left}}$. Line 404 decides whether the i th enqueue in b th block of internal node n is in the left child or right child subblocks of $n.\text{blocks}[b]$. By definition 8 to find an operation in subblocks of $n.\text{blocks}[i]$ we need to search in the range $n.\text{left}.\text{blocks}[n.\text{blocks}[i-1].\text{end}_{\text{left}}+1..n.\text{blocks}[i].\text{end}_{\text{left}}] \cup n.\text{right}.\text{blocks}[n.\text{blocks}[i-1].\text{end}_{\text{right}}+1..n.\text{blocks}[i].\text{end}_{\text{right}}]$.

There are $eb = n.\text{blocks}[b-1].\text{sum}_{\text{enq-left}}$ enqueues in the blocks 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 $i+eb$ 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 $n.\text{left}.\text{block}$ in the expected range of blocks, 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 left `Enqueues` are before the right `Enqueues`, there are $n.\text{blocks}[b].\text{num}_{\text{enq-left}}$ enqueues ahead from the left child. So we need to search for $i-n.\text{blocks}[b].\text{num}_{\text{enq-left}} + n.\text{blocks}[b-1].\text{sum}_{\text{enq-right}}$ (Line 409). Other parameters for the left child are chosen similar to the right child.

So, in both cases the direct subblock containing $E_i(n, b)$ is computed in Lines 405 and 409. Finally, $n.\text{child}.\text{GetENQ}(\text{subblock}, i)$ is invoked on the subblock containing $E_i(n, b)$ and it returns $E_i(n, b)$ by the hypothesis of the induction. \square

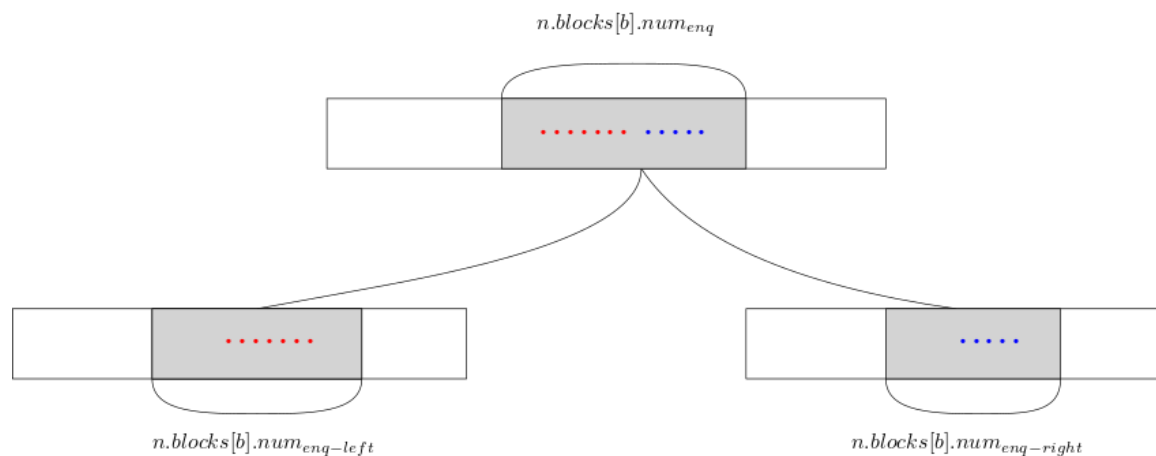


Figure 18: The number and ordering of the enqueue operations propagated from the left and the right child to $n.\text{blocks}[b]$. Enqueue operations from the left subblocks (colored red), are ordered before the enqueue operations from the right child (colored blue).

Definition 29. A Refresh is *successful* if it performs a successful CAS on Line 320. If Refresh instance R_1 does its CAS at Line 320 sooner than Refresh instance R_2 we say R_1 *happened* before R_2 .

Let i be the value R_n , a successful instance of Refresh() on the node n , reads from $n.head$. R_n does a successful CAS(null, new) into $n.blocks[h]$, where b is the new. Without loss of generality for the rest of this section assume n is the left child of $n.parent$. From now on we say p as an abbreviation for $n.parent$. Let R_p be the first successful $p.Refresh()$ that reads some value greater than i for $left.head$ and is contains b in its created block s in Line 317. From Lemma 24 we know there could be only one $p.Refresh()$ propagating b . R_p does a successful CAS(null, new) into $p.blocks[j]$, where s is the new.

Although other fields of b are set while creating it, because the index of the superblock of b is not known until it is propagated, R_n cannot set the `super` field of a b while creating it. One approach is to set the `super` field of b after it is propagated by R_b but this would not be efficient because there might be p subblocks in s . However, once b is installed, its superblock is going to be close to $n.parent.head$ 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. $b.super$ 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 $b.super$ to $n.parent.head$ while creating b , but the problem is that there might be many $p.Refreshes$ could happen that contain blocks from the right child of p and j could be arbitrarily (right word?) greater than $b.super$. We set $b.super$ to $p.head$ after appending b to $n.blocks$ (Line 326). Maybe R_n goes to sleep at some time after installing b and before setting $b.super$. In this case the next Refreshes on n and $n.parent$ help fill in the value of $b.super$.

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

Observation 30. After $n.blocks[i].CAS(null, b)$ succeeds, $n.head$ cannot increase from i to $i+1$ unless $b.super$ is set.

Proof. From the Observation 2 we know the only change to $n.head$ is on Line 327 which is incrementing. Before an instance of Advance() increments $n.head$ on Line 327, Line 326 ensures that $n.blocks[head].super$ was set at Line 326. □

Corollary 31. If $n.blocks[i].super$ is null, then $n.head < i$ and $n.blocks[i+1]$ is null.

Proof. If $b.super$ is null then $n.head$ cannot advance so the next $n.Refresh()$ will fail. By the previous corollary, $b.super$ has to be set before the next successful Refresh() on n after R_n . □

Now let us talk about how the $p.Refreshes$ that took place after the putting b into n , will help to set $b.super$ and propagate b .

Lemma 32. If $b \in n.parent.blocks[i]$ then $b.super \leq i$.

Proof. For R_p to contain block b , it has to read $n.head$ greater than h (see Line 334). For $n.head$ to be greater than h it means $n.head$ is incremented in Line 327 which means $b.super$ was already set in Line 326 (see Observation 30). So if R_p propagates b it means $b.super$ was already set. Let j be the value written in $b.super$. j has been read in Line 310 or Line 321 which both are before calling Advance that sets $b.super$. From Observation 2 we know $p.head$ is non-decreasing so $j \leq i$. The reader may wonder when the case $j = i$ happens, it happens when $p.blocks[j]=null$ while j is read and R_p puts its created block into $n.blocks[j]$. □

Lemma 33. If $R_n.Refresh()$ puts b in $n.blocks[h]$ at Line 320, then the block created by one of the next two successful $p.Refreshes$ according to the Definition 29 contains b and $b.super$ is set before Line 317 of the the second successful $p.Refresh()$.

Proof. It is sufficient to prove one of the two successful $p.Refresh()$ es propagate b . If the first successful $p.Refresh()$ propagated b then the claim is true, so in the remaining part we assume the first $p.Refresh()$ did not propagate b and prove the second $p.Refresh()$ propagates b .

$b.super$ is set by some instance of `Refresh()` on n or p showed by R' and $n.head$ is incremented by some `Refresh()` called R'' . We want to know how great $j - b.super$ can be. $p.head$ is hp when R' reads it. From Lemma 6 $p.head$ could only increase from hp to $hp+1$ if $p[hp] \neq \text{null}$. In other words there should be a successful $p.Refresh()$ for $p.head$ to increase. We claim there cannot be another successful $p.Refresh()$ after R' reads $p.head$ and before R_p performs Line 334.

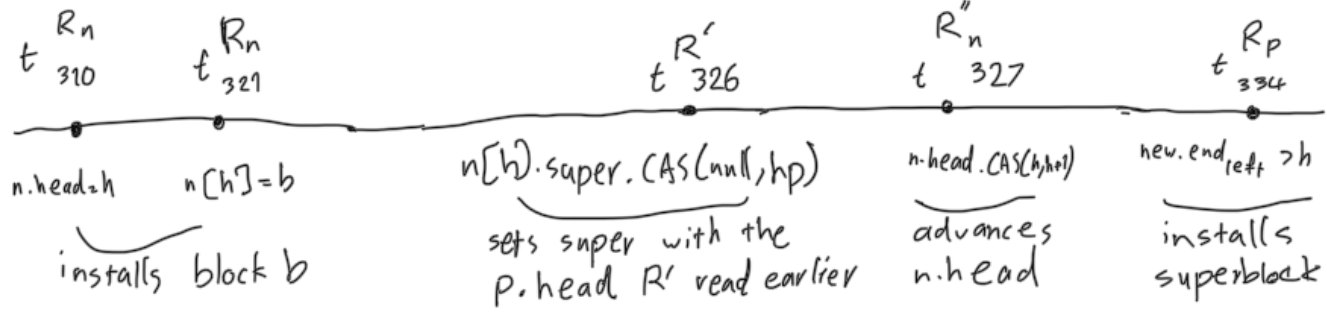


Figure 19: Time relations between R_n, R_p, R', R''

Assume the first successful $p.Refresh()$ after $t_{320}^{R_n}$ did not set $b.super$. It might happen maybe because the value read for h_{left} in Line 312 is less than i or maybe $i = h_{left}$ and $left.blocks[h_{left}] = \text{null}$, which means $n.head$ is advanced but b is still not installed in $n.blocks[i]$ which means R_n has not reached to the Line 320.

Let the first successful $p.Refresh()$ be R_{p1} and the second next successful $p.Refresh()$ be R_{p2} . If R_{p1} reads x in Line 310, then R_{p2} has to read $x+1$ in Line 310 (induced from 6, 2). See the timeline in Figure 20 for two consecutive successful `Refresh()` instances R_{p1} , R_{p2} on p .

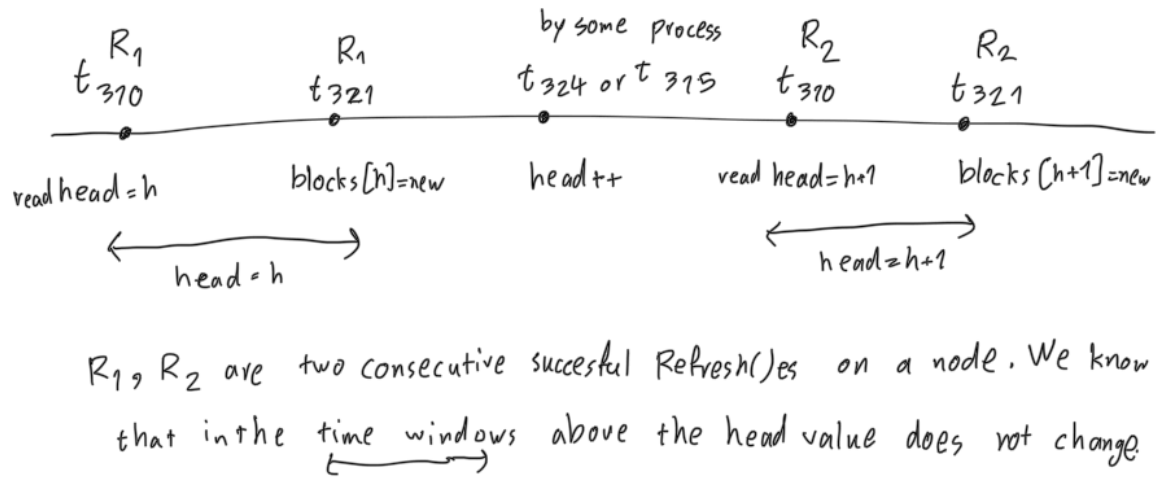


Figure 20: R_{p2} reads $p.head$ after $t_{321}^{R_{p1}}$, which is after $t_{321}^{R_n}$. R_{p2} has to help increment $n.head$ and set $b.super$.

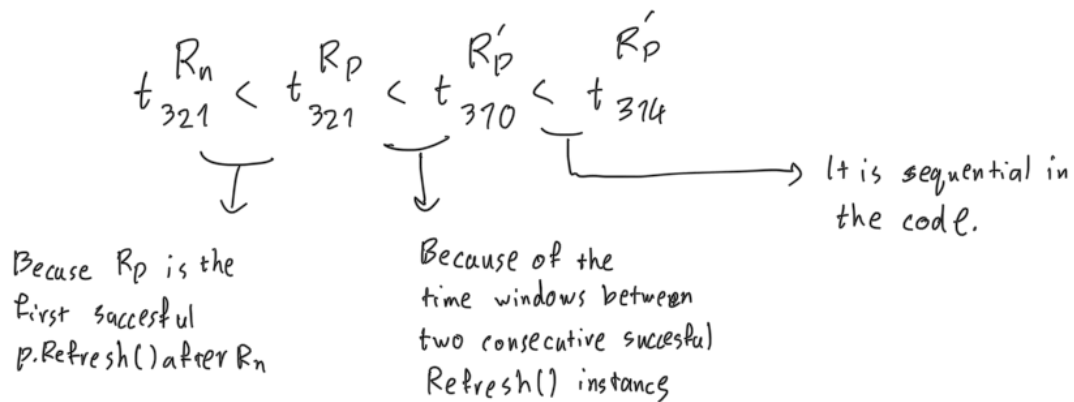


Figure 21: The second `Refresh` on p contains b and reads $n.head > i$.

So `b.super` has set by some process before the second next successful `p.Refresh()` on Line 326. Since `i` is read in the Line 310 then the `CreateBlock()` in Line 317 is going to read some value for `left.head` greater than `h` and propagates `b` to `p`. So if `b` was not propagated already we are sure the second next successful `p.Refresh()` propagates `b`. \square

Corollary 34. *If `n.blocks[b]` is propagated to `n.parent`, then `n.blocks[b].super` is equal to or one less than the index of the superblock of `n.blocks[b]`.*

Now using Corollary 34 on each step of the `IndexDEQ` we prove its correctness.

Lemma 35 (`IndexDEQ` correctness). *If $1 \leq i \leq n.blocks[b].num_{deq}$ then `n.IndexDEQ(b, i)` returns the rank in $D(root)$ of $D_i(n, b)$.*

Proof. We will prove this by induction on the distance of n from the `root`. The base case where n is `root` is trivial (Line 415). For the non-root nodes `n.IndexDEQ(b, i)` computes the superblock of the b th block of n in Line 418 by Corollary 34. So preconditions of the invocation of `BinarySearch` is satisfied. After that the order of $D_i(n, b)$ in $D(n.parent, superblock)$ is computed in Lines 419–425. Note that by Lemma 25 in each block there is at most one `Dequeue` from each process and by Definition 15 `Dequeues` in a block are ordered based on the order of its subblocks from left to right. Then 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 422). Finally `IndexDEQ` is called on `n.parent` recursively and it returns the response from induction hypothesis. \square

We have explained the picture of proof in the first paragraph of proof section. Should I repeat here?

Lemma 36 (satisfiability). *L can be a linearization ordering.*

Proof. A linearization for an execution is legal if, op_1 terminates before op_2 starts in execution then op_1 is before op_2 in the linearization. If op_1 terminates before op_2 starts, $op_1.\text{Append}()$ has terminated before $op_2.\text{Append}()$ started. From Lemma ?? 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 linearized in any order among themselves. The way we have chosen to order the operations in a block makes `GetENQ`, `FindResponse` efficient.

Definition 37. If a `Dequeue` operation returns null it is called *null dequeue*, otherwise it is called *non-null dequeue*.

Observation 38. *In a sequential execution on a queue, kth non-null dequeue returns the element of kth enqueue.*

Definition 39. Assume the operations in `root.blocks` are applied on an empty queue in order of L . If element of `Enqueue` e is the response to `Dequeue` d , we say $\text{Resp}(d) = e$, otherwise if d is null dequeue then $\text{Resp}(d) = \text{null}$.

Observation 40. *By Definition 15 enqueue operations come before dequeue operations in a block.*

In the next Lemma we talk about what the `size` field in `root` blocks store.

Lemma 41. `root.blocks[b].size` is the size of the queue if the operations in `root.blocks[0...b]` are applied in order of L .

Proof. We prove the claim by induction on b . Base case $b = 0$ is trivial since the queue is initially empty and `root.blocks[0].size` = 0. We are going to show the correctness of case $b = i$ assuming the correctness of $b = i - 1$. By 40 If there are more than `root.blocks[i-1].size + root.blocks[i].sum_enq` dequeue operations in `root.blocks[i]` then the queue would become empty after `root.blocks[i]`, otherwise size of the queue after b th block in root is `root.blocks[b-1].size + root.blocks[b].sum_enq - root.blocks[b].sum_deq` which is same as the assignment on Line 344. \square

Next Lemma is about the relation between `#enqueues`, `#non-null dequeues` and `size`.

Lemma 42. *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].sum_enq - root.blocks[b].size`*

Proof. There are `root.blocks[b].sum_enq` enqueue operations in `root.blocks[0...b]`. Size of the queue after doing `root.blocks[0...b]` in order L is `#enqueues - #non-null dequeues` in `root.blocks[0...b]`. By correctness of `size` field from Lemma 41 `#non-null dequeues` is `root.blocks[b].sum_enq - root.blocks[b].size`. \square

Corollary 43. *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].num_enq - root.blocks[b].size + root.blocks[b-1].size`*

Lemma 44. $R(D_{\text{root},b,i})$ is null iff `root.blocks[b-1].size + root.blocks[b].num_enq - i < 0`.

Proof. From 43, 37 and 40. \square

Lemma 45 (Computing Response). `FindResponse(b, i)` returns $\text{Resp}(D_i(\text{root}, b)).\text{element}$.

Proof. Note that by Definition 15 the linearization ordering of operations will not change as new operations come so instead of talking about the linearization of operations before the $E_i(\text{root}, b)$ we talk about what if the whole operation in the linearization are applied on a queue.

$D_i(\text{root}, b)$ is $D_{\text{root.blocks}[b-1].\text{sum_deq}+i}(\text{root})$ from the Definitions 15 and sum_deq . $D_i(\text{root}, b)$ returns null if `root.blocks[b-1].size + root.blocks[b].num_enq - i < 0` by Lemma 44 (Line 220). Otherwise if it is k th non-null dequeue in L it returns k th enqueue by

Observation 38. By Lemma 42 there are $\text{root.blocks}[b-1].\text{sum}_{\text{enq}} - \text{root.blocks}[b-1].\text{size}$ non-null dequeue operations in $\text{root.blocks}[0 \dots b-1]$. Note that the dequeues in $\text{root.blocks}[b]$ before $D_i(\text{root}, b)$ are non-null dequeues. So the response is $E_{i-\text{root.blocks}[b-1].\text{size}+\text{root.blocks}[b-1].\text{sum}_{\text{deq}}}(\text{root})$ (Line 222). See figure 22.

After computing e , the code finds b, i such that $E_i(\text{root}, b) = E_e(\text{root})$ by `DSearch` and then finds its element using `GetEnq` (Line 223). □

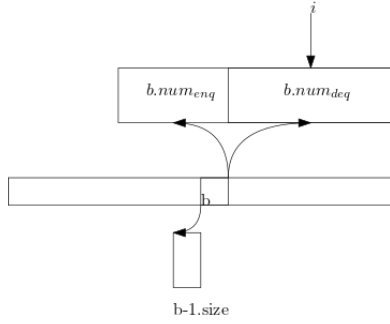


Figure 22: The position of $D_i(\text{root}, b)$.

Lemma 46 (correctness). *If operations are applied as L on a sequential queue, the sequence of the responses would be the same as our algorithm.*

Proof. Enqueue operations do not return any value so they are fine. Response of a dequeue in sequential execution of L is same as our algorithm from 45. □

Theorem 47 (Main). *The queue implementation is linearizable.*

Proof. From 36, 46. □

Lemma 48 (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 dequeue in $D(n)$ there are some constant steps and IndexDeq() is called. IndexDeq does a search with range p in each level (Lemma ??) 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 27, 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 GetEnq() takes $O(\log p)$ since there are $\leq p$ subblocks in a block (Lemma 26), so GetEnq() 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 49 (CASes invoked). *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() has 2 CASes, one in Line 320 and one in Lines ?? or ??. \square

Lemma 50 (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 $\text{end} - b - 1 \leq 2 \times \text{root.blocks}[b].\text{size} + \text{root.blocks}[\text{end}].\text{size}$. 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 $\text{root.blocks}[\text{end}]$'s first Dequeue(). Furthermore in the execution of queue operations in the linearization ordering, the size of the queue becomes $\text{root.blocks}[\text{end}].\text{size}$ after the operations of $\text{root.blocks}[\text{end}]$. There can be at most $\text{root.blocks}[b].\text{size}$ Dequeues in $\text{root.blocks}[b + 1 \dots \text{end} - 1]$; otherwise all elements enqueued by $\text{root.blocks}[b]$ would be dequeued before $\text{root.blocks}[\text{end}]$. The final size of the queue after $\text{root.blocks}[1 \dots \text{end}]$ is $\text{root.blocks}[\text{end}].\text{size}$. After an execution on a queue the size of the queue is greater than or equal to $\# \text{enqueues} - \# \text{dequeues}$ in the execution. We know the number of dequeues in $\text{root.blocks}[b + 1 \dots \text{end} - 1]$ is less than $\text{root.blocks}[b].\text{size}$, therefore there cannot be more than $\text{root.blocks}[b].\text{size} + \text{root.blocks}[\text{end}].\text{size}$ Enqueues in $\text{root.blocks}[b + 1 \dots \text{end} - 1]$. Overall there can be at most $2 \times \text{root.blocks}[b].\text{size} + \text{root.blocks}[\text{end}].\text{size}$ operations in $\text{root.blocks}[b + 1 \dots \text{end}]$ 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 $2 \times \text{root.blocks}[b].\text{size} + \text{root.blocks}[\text{end}].\text{size}$ blocks in between $\text{root.blocks}[b]$ and $\text{root.blocks}[\text{end}]$. So $\text{end} - b - 1 \leq 2 \times \text{root.blocks}[b].\text{size} + \text{root.blocks}[\text{end}].\text{size}$.

So the doubling search reaches start such that the $\text{root.blocks}[\text{start}].\text{sum}_{\text{enq}}$ is less than e in $O(\log(\text{root.blocks}[b].\text{size} + \text{root.blocks}[\text{end}].\text{size}))$ steps. See Figure 23. After Line 805, the binary search that finds b also takes $O(\log(\text{root.blocks}[b].\text{size} + \text{root.blocks}[\text{end}].\text{size}))$. Next, i is computed via the definition of sum_{enq} in constant time (Line 807). So the claim is proved. \square

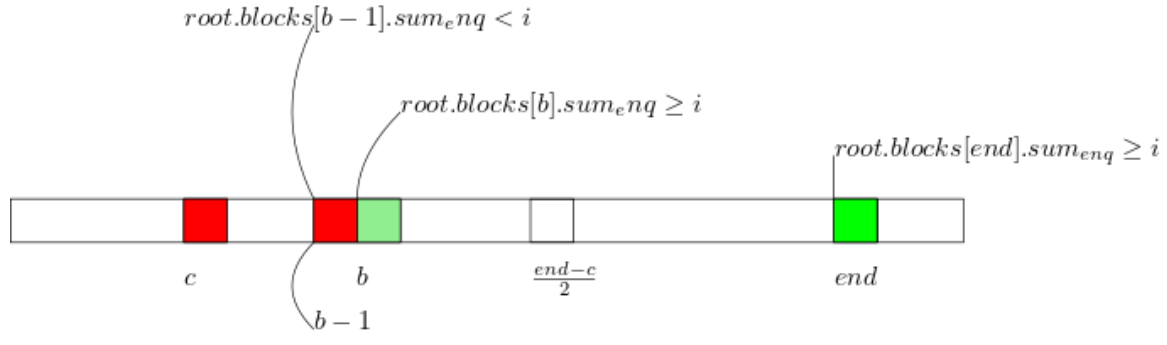


Figure 23: Distance relations between `start`, `b`, `end`. `start` is shown with `c`.

4.1 Garbage Collection or Getting rid of the infinite Arrays

5 Using Queues to Implement Vectors

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

6 Conclusion

possible directions for work

Maybe Stacks

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