

## 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.

## 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.

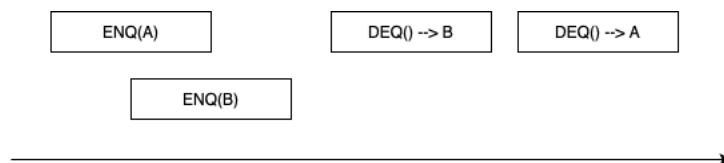


Figure 1: An example of a linearizable execution. Either `Enqueue(A)` or `Enqueue(B)` could take effect first since they are concurrent.

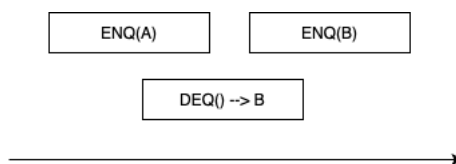


Figure 2: An example of an execution that is not linearizable. Since `Enqueue(A)` has completed before `Enqueue(B)` is invoked the `Dequeue()` should return A or nothing.

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.

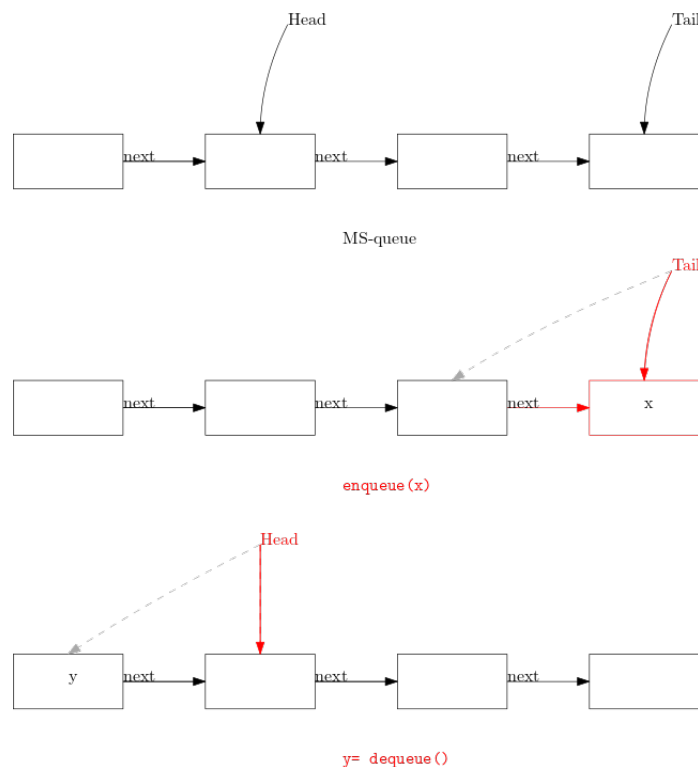


Figure 3: MS-queue structure, enqueue and dequeue operations. In the first diagram the first element has been dequeued. Red arrows show new pointers and gray dashed arrows show the old pointers.

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.

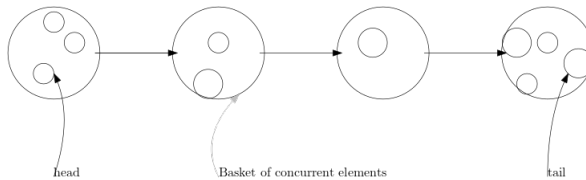


Figure 4: Baskets queue idea. There is a time that all operations in a basket were running concurrently, but only one has succeeded to do CAS. To order the operations in a basket, the mechanism in the algorithm for processes is to CAS again. The successful process will be the next one in the basket and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

Ladan-Mozes and Shavit [?] 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 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  $\Omega(p)$ .

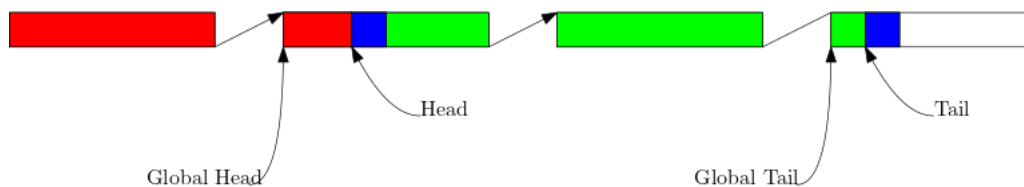


Figure 5: Global pointers point to arrays. Head and Tail elements are blue, dequeued elements are red and current elements of the queue are green.

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.

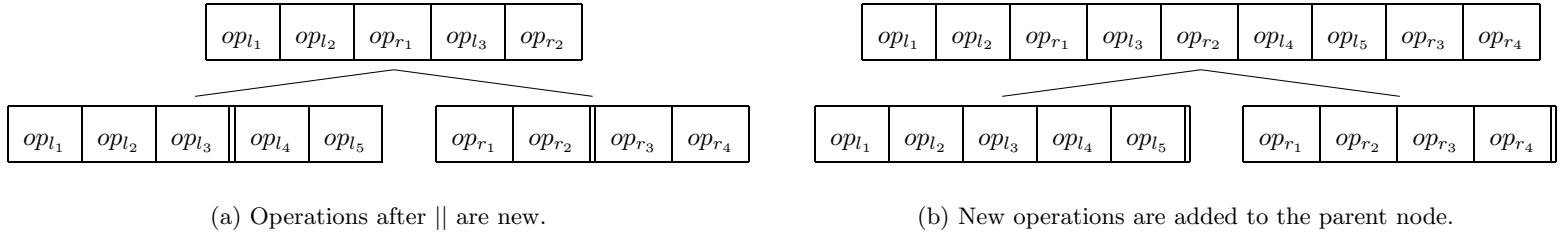
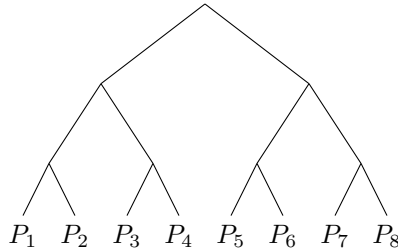


Figure 6: Propagate Step in Universal Construction

### 3 Universal Construction using Tournament Tree with Big CAS Objects

A universal construction gets any sequential object and creates a concurrent version of the given datastructure. Our universal construction in Algorithm 1 relies on a tournament tree with  $p$  leaves and height  $\log(p)$  is shared among  $p$  processes. Nodes in the tree are CAS objects that store an ordering and support **append()** and **diff()** operations. Leaf  $l_i$  stores the sequence of the operations invoked by  $P_i$ . Each internal node stores the sequence of operations propagated up to it. When process  $P_i$  wishes to apply an operation  $op$  on the implemented object, it appends  $op$  to its assigned leaf and tries to propagate it up to the root. The history of operations stored at the root is the linearization ordering. The operation  $op$  is linearized when it is appended to the root.



The algorithm uses a subroutine **REFRESH**( $n$ ) that concatenates new operations from node  $n$ 's children (that have not already been propagated to  $n$ ) to the sequence of operations stored in  $n$  and tries to CAS the new sequence into  $n$ . In other words, **REFRESH**( $n$ ) tries to append  $n$ 's children's new operations to  $n$ 's sequence. After a process adding a new operations to its leaf, it has to propagate new operations up to the root. **PROPAGATE**( $n$ ) tries to append  $n$ 's new operations to the root  $n$  by recursively calling **REFRESH**( $n$ ). In each **Propagate**() step if a **REFRESH**( $n$ ) fails, it means another CAS operation has succeeded; if so, it tries to **REFRESH**( $n$ ) again. If the second attempt fails too, another process has already appended the operations the current **PROPAGATE** is trying to append. Operations that were in  $l_i$  before **PROPAGATE**( $l_i$ .parent) was invoked are guaranteed to be added to the root by the time the **PROPAGATE**( $l_i$ .parent) terminates.

---

#### Algorithm Universal Construction Idea

---

1: <i>response</i> Do(operation $op$ , $pid$ $i$ )	14: <i>boolean</i> REFRESH( <i>node</i> $n$ )
2: $l_i$ .APPEND( $op$ )	15: $old = READ(n)$
3:   PROPAGATE(parent of $l_i$ )	16: $new =$ ops that $n$ 's children contain but $old$ does not
4:   Run the sequence stored in root	17: $new = old \cdot new$
5: <b>return</b> $op$ 's response from line 4	18: <b>return</b> $n.CAS(old, new)$
6: <b>end</b> Do	19: <b>end</b> REFRESH
7: <i>void</i> PROPAGATE( <i>node</i> $n$ )	
8: <b>if</b> $n == root$ <b>then return</b>	
9: <b>else if</b> !REFRESH( $n$ ) <b>then</b>	
10:     REFRESH( $n$ )	
11: <b>end if</b>	
12:   PROPAGATE(parent of $n$ )	
13: <b>end</b> PROPAGATE	

---

$O(\log n)$  CAS operations are invoked to do a PROPAGATE, but the CAS words store sequences of unbounded length. The problem is that we are trying to store unbounded sequence of operations in each node  $n$  (see Figure 7). However, to compute the result of an operation, we only use the total ordering that is stored at the root. Although we use a similar construction for our queue implementation, we develop an implicit representation of the sequence of operations, so that we can use reasonable sized CAS objects and still achieve polyarithmic step complexity.

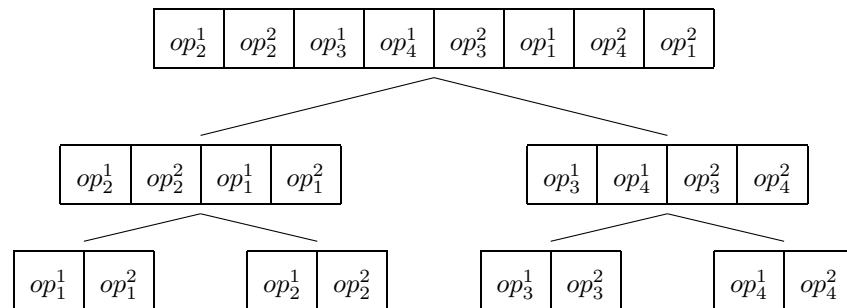


Figure 7: Universal Construction:  $op_j^i$  denotes the  $i$ th operation from process  $j$ . In each node, we store the ordering of all the operations propagated up to it.

## 4 Block Trees

Here we introduce a data structure that allows processes to agree on the linearization ordering of their operations using  $O(\log p)$  CAS per operation called a *block tree*. Then we use the block tree as a stepping stone towards our queue algorithm. A block tree is a tournament tree shared among  $p$  processes (see Figure 8). Each process has a leaf, and it appends its operations to its leaf. After that, the process tries to propagate its new operation up to the tree's root. An ordering of operations propagated up to a node is stored in that node. All processes agree on the sequence stored in the root and this is used as the linearization ordering. Our idea is similar to Jayanti and Petrovic's multi-enqueuer single-dequeuer Queue [?], but we do not use CAS operations with big words and do not put a limit on the number of concurrent operations.

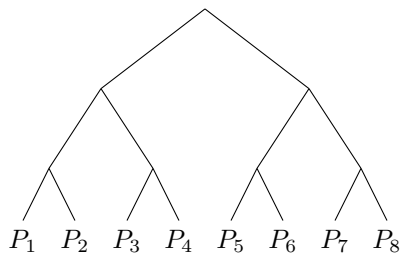


Figure 8: In the block tree each process has a leaf and in each node there is an ordering of operations stored. Each node tries to propagate its operations up to the root, which stores the final ordering of all operations.

The goal here is to ensure that in each propagate step the new operations are propagated up to the parent in  $O(\log p)$  steps (see Figure 9). Then, a dequeue operation uses the linearization ordering to compute its answer.

In each propagate step, our algorithm uses a subroutine  $\text{REFRESH}(n)$  that aggregates new operations from node  $n$ 's children (that have not already been propagated to  $n$ ) and tries to append them into  $n$  using a CAS operation. The general idea is that if we call  $\text{REFRESH}(n)$  twice, the operations in  $n$ 's children before the first  $\text{REFRESH}(n)$  are guaranteed to be in  $n$ . Instead of storing operations explicitly in the nodes, we only keep track of the number 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  $\text{Refresh}()$  procedure and only propagate the block of the new operations.

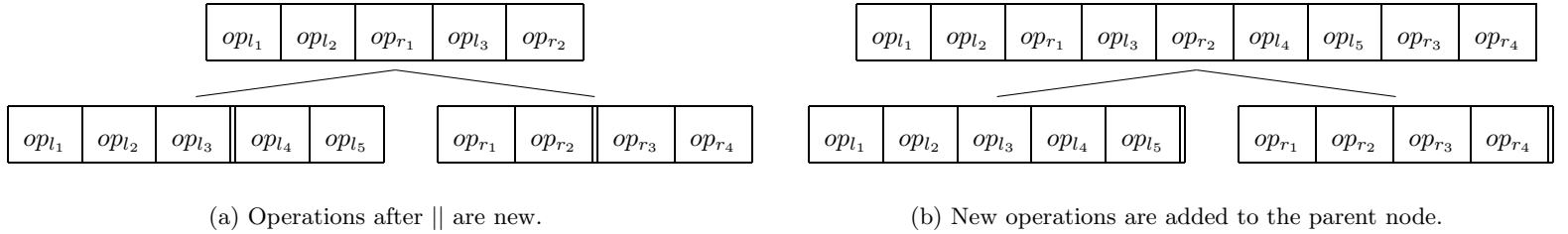


Figure 9: Successful **Refresh**, operations in children after || are new.

We also implement methods **Get(i)**, **Index(op)** to get the  $i$ th propagated operation and compute the rank of a propagated operation in the linearization. **Get(i)** finds the block containing the  $i$ th operation in the root and then 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)**.

**Get()** and **Index()** search among blocks in each level of the tree to find the sub-block or super-block containing the given operation. Each block stores a constant amount of information (like prefix sums) to allow binary searches to find the required block in a node quickly.

Block tree can be used to implement queue, but **Get(i)** may take a long time since it has to find the block containing the  $i$ th operation at the root level.

We apply two ideas from universal construction to create a new linearizable data structure agreeing on a sequence of elements among processes. First, there is a shared tournament tree among processes, in which each process appends its element to its leaf in the tree and then tries to propagate it up to the root by performing **REFRESH()** operations at each node. Second, each operation is linearized when its element is appended to the root.

In the universal construction, we order new concurrent operations at each **REFRESH()** and maintain that order in the path up to the root. However, we can instead keep track of sets of concurrent operations and create the total ordering of all operations at the root (see Figure 10).

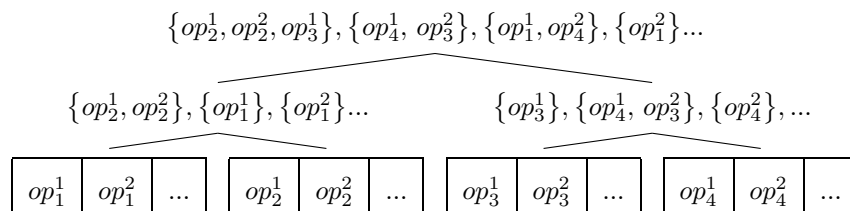


Figure 10: In each internal node, we store the set of all the operations propagated together, and one can arbitrarily linearize the sets of concurrent operations among themselves. Since we linearize operations when they are added to the root, ordering the blocks in the root is important.

The definition of linearizability allows concurrent operations to be reordered arbitrarily. Thus, a group of concurrent operations can be appended to our root sequence as one block without specifying the order among the operations.

We used unbounded CAS objects storing sequences as big words in the universal construction. One can represent sequences as arrays to overcome this implementation problem. Each array element will store one of the blocks of concurrent operations described in section ??.

Copying operation sequences from children to their parent in a **REFRESH()** takes time proportioned to the number of operations being copied. This is time-consuming, so we propose a way to augment the tree to calculate lines 15,16 in  $O(\log p)$  steps which reads new



operation and concatenates them with old operations. Instead of representing the set of operations by explicitly listing them in a node, we represent a set of operations implicitly by recording which of the children's sets were unioned to create the set. Having operation sequences stored at leaves, we can deduce a set of operations in a node using this implicit representation. (see Figure 11.)

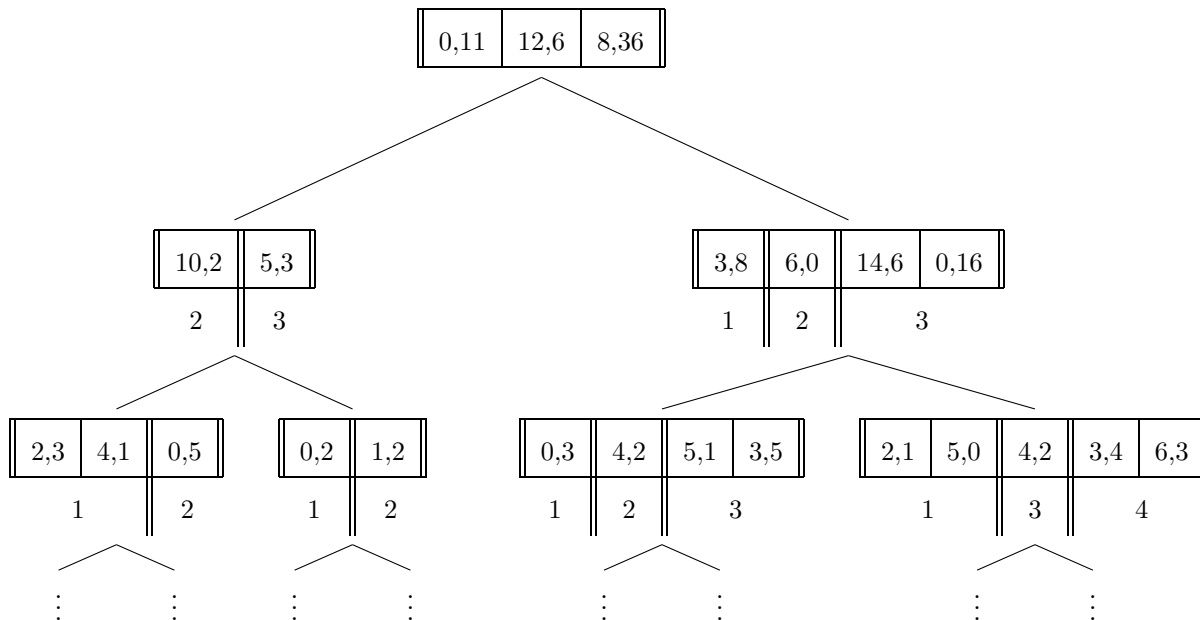


Figure 11: Showing concurrent operation sets with blocks. Each block consists of a pair(left, right) indicating the number of operations from the left and the right child, respectively. Block (12,6) in the root contains blocks (10,2) from the left child and (6,0) from the right child. Blocks between 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).

**Definition 1.** {Existence of an operation in a block} Operation  $op$  exists in block  $b$  if it has propagated up to block  $b$ .

**Definition 2.** {Subblock} The blocks that are aggregated into block  $b$  in a `PROPAGATE()` step are called subblocks of  $b$ . Block  $b_1$  is a subblock of  $b_2$  if and only if  $b_1$  is a block in node  $v$  and in  $b_2$  is a block in the parent of  $v$  and  $b_1$ 's elements exists in  $b_2$ 's elements.

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 with come before the right child and the order of blocks of each child is preserved among themselves.

In a `PROPAGATE()` invocation path from a leaf to root, there will be `REFRESH()` steps with merges from  $2, 4, 8, \dots, p$  processes. So in a complete propagation, at most  $2p$  blocks are merged into one block. (maybe useful for analysis)

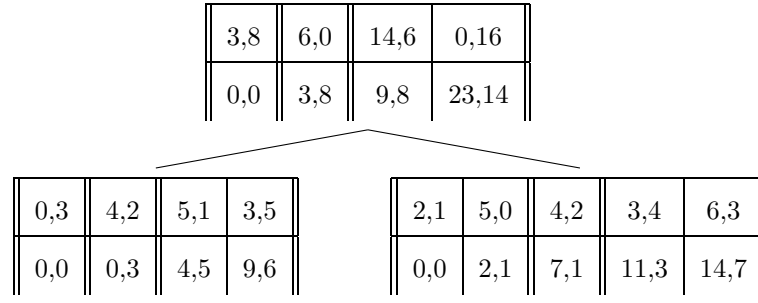


Figure 12: Using Prefix sums in blocks. When we want to find block  $b$  elements in its children, we can use binary search. The number below each block shows the count of elements in the previous blocks.

#### 4.1 Using pointers and prefix sum to make $\text{GetIndex}(i)$ faster

$\text{GETINDEX}(i)$  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 BSearch 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 13 and 12).

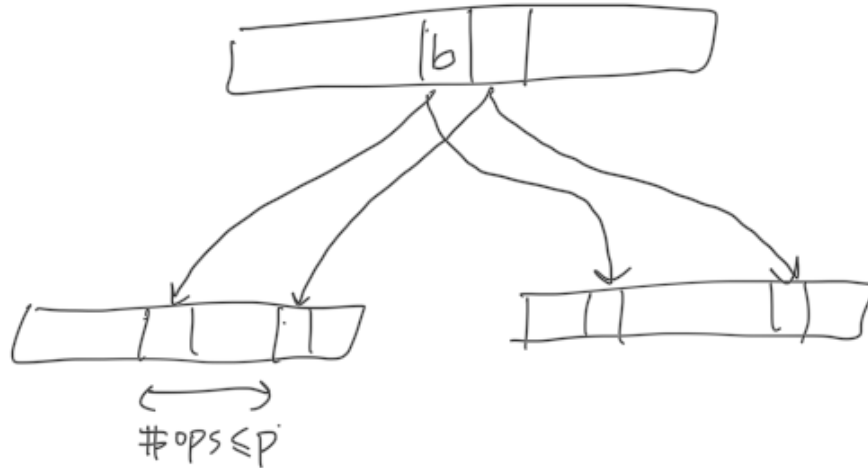


Figure 13: Block have pointers to the starting block of theirs for each child.

Starting from the root,  $\text{GETINDEX}(i)$  BSearches  $i$  in the prefix sum array to find block containing  $i$ th operation, then continues recursively calling  $\text{GETELEMENT}(b, i)$  to find  $i$ th element of block  $b$ . From lemma ?? we know a block size is at most  $p$ . So BSearch 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  $O(p)$ .

## 4.2 Block Tree Algorithm

Our Block Tree is a linearizable implementation of a data structure that stores a sequence of elements. It has two methods (see Algorithm ??), `APPEND( $e$ )` which appends element  $e$  to the sequence, and `GET( $i$ )` which returns the  $i$ th element in the sequence.

**Design of a block tree** Each process is assigned to a leaf in a shared tournament tree. Thus, for example, the leaf node for process  $p_i$  contains an array of elements by  $p_i$  in the order they were invoked. Each internal node of the tree contains an array of blocks of elements. Block  $b$  in node  $n$  is created in a `PROPAGATE()` step and is merged block of new blocks at the time of `PROPAGATE()` reading  $n$ 's children blocks. Each block consists of pointers left and right, to the last block merged into itself from left and right child in that order. Moreover, two numbers, left and right, indicate the count of elements in the blocks from the left and right child consecutively. Furthermore, prefix left, and right can be computed from the prefix sum of left and right values. Elements of block  $b$  can be determined recursively (`GETELEMENTS( $b$ )`). The  $i$ th element in the sequence can be determined in  $O(\log^2 p)$  steps by recursively finding  $i$ th element in block  $b$  (`GETELEMENT( $i$ )`). After element  $e$  is propagated (appended to a block into the root), its index can be computed with `GETINDEX( $op$ )`.

In order to compute elements of a block faster we store prefix-sum blocks (block  $i$  has tuple(right-sum=#right ops in previous block, left-sum=#left ops in previous blocks)[See Figure 12]. Here is the algorithm to get elements of a block.

**Specification** A block tree is a shared data structure that stores a sequence of elements. It has two methods `Append( $e$ )` and `Get( $i$ )`. `Append( $e$ )` adds  $e$  to the end of the sequence and returns the index of  $e$  in the sequence. `Get( $i$ )` returns  $i$ th element stored in the sequence.

**SubBlock** Block  $s$  is a subblock of  $b$  if  $s$  is between blocks `start..end` in  $n$  from Lines 41,42 of `CreateBlock()`.

**Membership** Element  $e$  is a member of block  $b$  in:

- internal node  $n$ , if  $e$  is a member of  $s$  that  $s$  is a subblock of  $b$ .
- leaf node  $n$ , if  $e$  belongs to  $n.dir.blocks[b'.end_{dir}+1..b.end_{dir}]$  for  $dir \in \{left, right\}$  which  $b'$  is the previous block of  $b$  in  $n$ .

**Order of elements inside node** Element **d** is before element **e** in node **n**, if:

- The block containing **d** is before the block containing **e**.
- **e** and **d** are in the same block and **d** is in the left child and **e** is in the right child.
- **d** is before **e** in the same child's order.

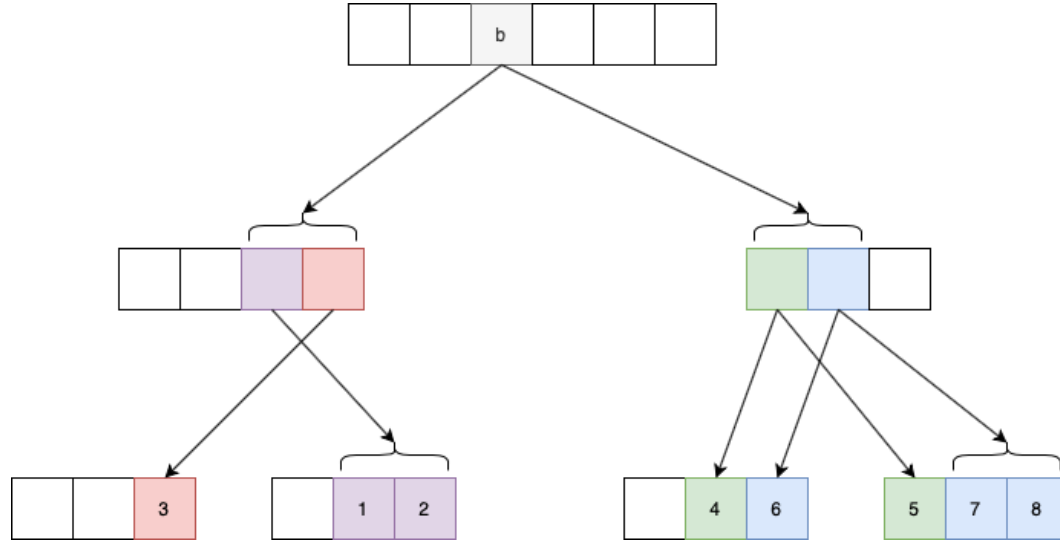


Figure 14: Order of elements in **b**: elements in leaves are ordered with numerical order in the drawing.

**CreateBlock()** **CreateBlock(n)** returns a block containing new operations of **n**'s children. **b'.end<sub>left</sub>** stores the index of the rightmost subblock of left child of **b**'s previous block. Other attributes are assigned values followed by definition.

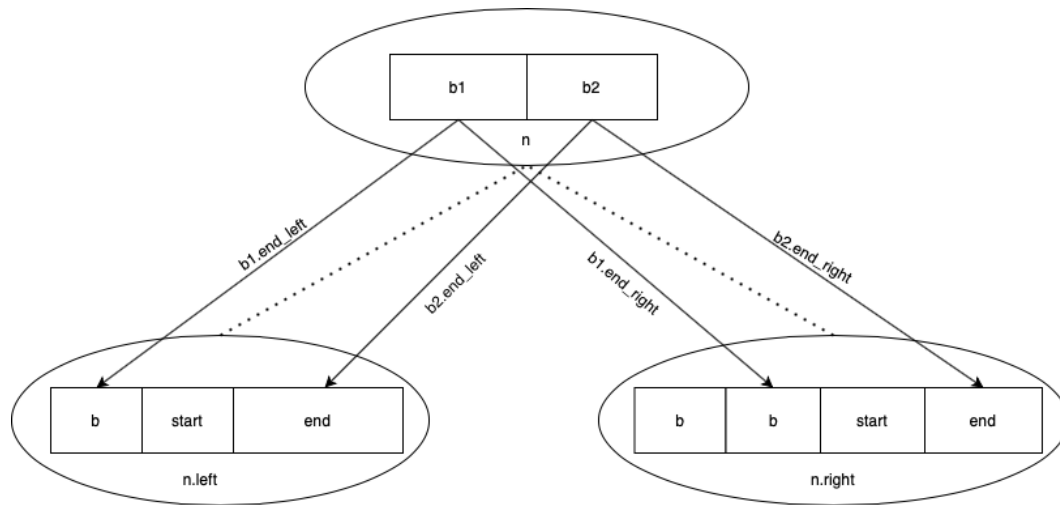


Figure 15: Snapshot of a **CreateBlock()**

**Double Refresh** Elements in  $n$ 's children's blocks before line 13 are guaranteed to be in  $n$ 's blocks after Line 15.

*Proof.* `CreateBlock()` reads blocks in the children that does not exist in the parent and aggregates them into one block. If a `Refresh` procedure returns true it means it has appended the block created by `CreateBlock()` into the parent node's sequence. So suppose two `Refreshes` fail. Since the first `Refresh` was not successful, it means another CAS operation by a `Refresh`, concurrent to the first `Refresh`, was successful before the second `Refresh`. So it means the second failed `Refresh` is concurrent with a successful `Refresh` that assuredly has read block before the mentioned line 13. After all it means if any of the `Refresh` attempts were successful the claim is true, and also if both fail the mentioned claim still holds.  $\square$

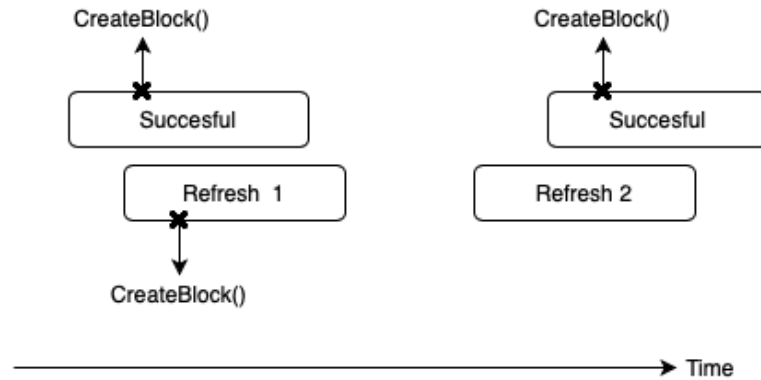


Figure 16: The second failed `Refresh` is assuredly concurrent to a Successful `Refresh()` with `CreateBlock` line after first failed `Refresh`'s `CreateBlock()`.

**Disjunction** Blocks in node  $n$ 's contain disjoint sets of elements.

*Proof.* Without loss of generality, assume blocks  $b_1$ ,  $b_2$  contain common element  $e$  from the left child, and  $b_2$  is after  $b_1$  in  $n$ 's sequence of blocks. So block start of  $b_2$ 's `CreateBlock()` is after block end of  $b_1$ 's end. Since  $b_2$ 's start is the end of the block before itself, it cannot be before  $b_1$ 's end.  $\square$

**Total Order** Sequence represented by the Block Tree is the sequence of the blocks stored in the root.

**Linearization Points** `Get(i)` is linearized when it terminates. `Append(e)` is linearized right after when a block containing  $e$  is appended to the root, if there are multiple elements appended together, they are linearized by the defined order in the root.

**Subblocks Upperbound** Block  $b$  has at most  $p$  subblocks.

*Proof.* If there are more than  $p$  subblocks, then there is more than one block from process  $p_l$ . `Append(e)` finishes after propagating and appending  $e$  to the root(line 9). So these blocks cannot be appended to root already, so  $p_l$  has invoked two concurrent `Append()`s(line 1) without terminating the first one.  $\square$

**Computing `Get(n, b, i)`** To find the  $i$ th element in block  $b$  of node  $n$ , we search among subblocks of  $b$  that is bounded by  $p$ . Subblocks of a block are within the start and end block of the `CreateBlock()` procedure of it.

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

1. Read  $n$ 's counter and head
2. Create block  $b$
3. CAS  $b$  into  $n$

4. If previous succeed:

- (a) Update sup of b's ending subblocks
- (b) Increment children's counters

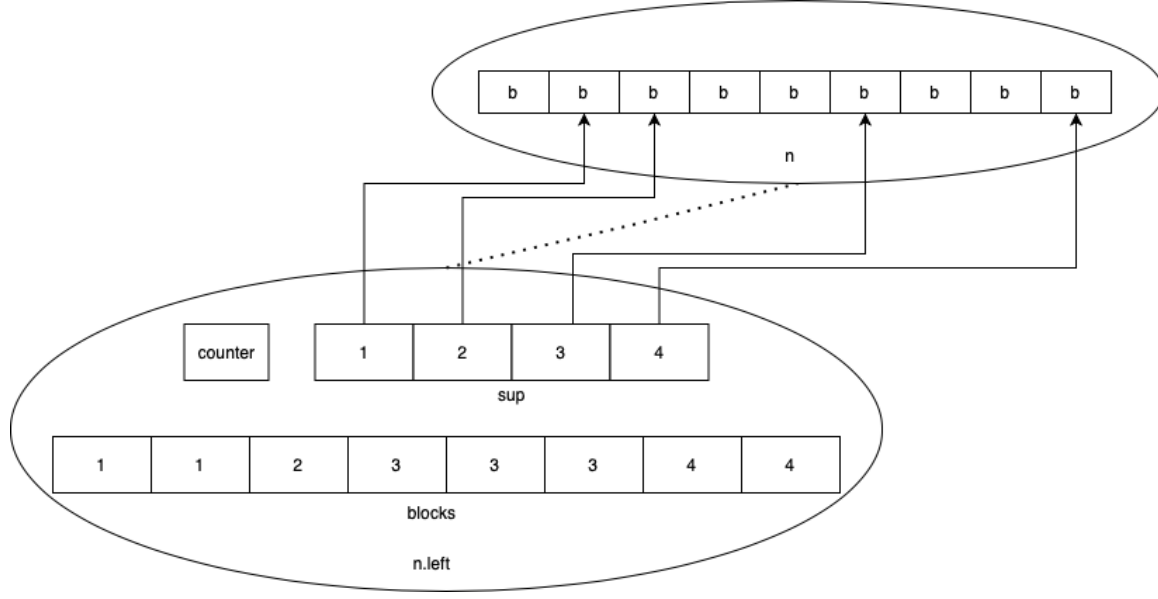


Figure 17: Sup and timer in a node, numbers on blocks are their time values.

### Computing superblock

1. Value read for `super[b.time]` in line 71 is not null.

*Proof.* `Index()` is invoked after finishing `Propagate()` in line 10. For each value `c_dir` read in lines 23, `super` is set before incrementing in lines 26,27. □

2. `super[]` preserves order from child to parent; if in a child block `b` is before `c` then `b.time ≤ c.time` and `super[b.time] ≤ super[c.time]`

*Proof.* Follows from the order of lines 37, 26, 27. □

3. `super[i+1]-super[i] ≤ p`

*Proof.* In a Refresh with successful CAS in line 24, `super` and `counter` are set for each child in lines 26,27. Assume the current value of the counter in node `n` is `i+1` and still `super[i+1]` is not set. If an instance of successful `Refresh(n)` finishes `super[i+1]` is set a new value and a block is added after `n.parent[sup[i]]`. There could be at most  $p$  successful unfinished concurrent instances of `Refresh()` that have not reached line 27. So the distance between `super[i+1]` and `super[i]` is less than  $p$ . □

4. Superblock of `b` is within range  $\pm 2p$  of the `super[b.time]`.

*Proof.* `super[i]` is the index of the superblock of a block containing block `b`. It is trivial to see that `n.super` and `n.b.counter` are increasing. `super(b)` is the real superblock of `b`. `super(t)` is the index of the superblock of the last block with time `t`. If `b.time` is `t` we have:

$$super[t] - p \leq super[t-1] \leq super(t-1) \leq super(b) \leq super(t+1) \leq super(t+1) \leq super[t] + p$$

□

## 5 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.

Next, we describe how to use block tree to implement queues. The block tree, maintains the history of all operations, not only the current state of the queue. Now consider the following history of operations. What should each **Dequeue()** return? We can implement Enqueue and Dequeue using our block tree. An **Enqueue(e)** appends an operation with input argument **e** in the block tree. To do 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

**Definition 3.** 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 ?? 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 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:

`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 dequeus[b-1]+ index of DEQ in the block's dequeus]`  
`: null;`

## 6 Main Algorithm

**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.

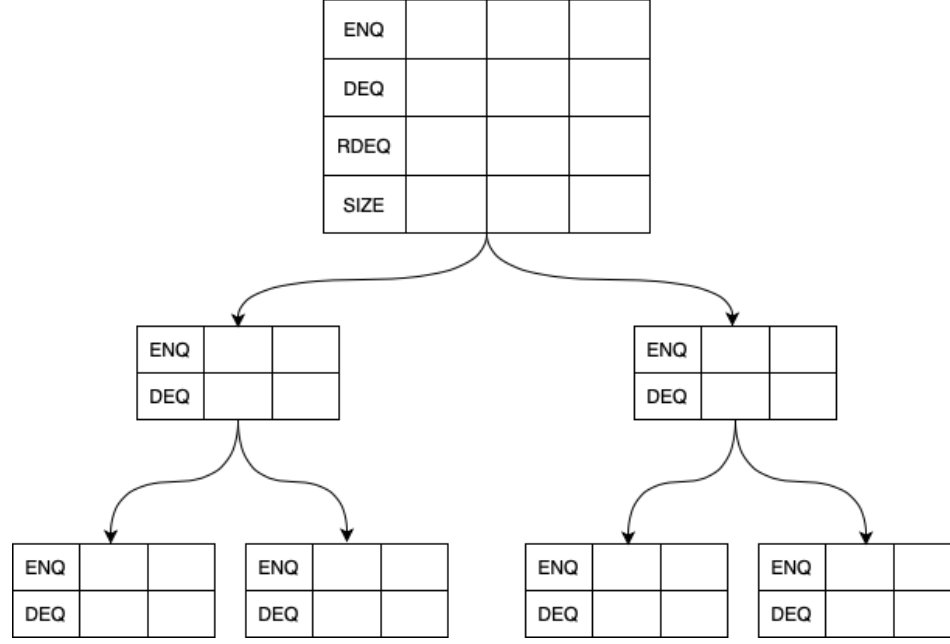


Figure 18: Fields stored in the Queue nodes.

### 6.1 Pseudocode description

**Tournament 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.

**Definition 4** (Block). A block is an object that stores some statistics described in Algorithm Queue.

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

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 15 .

**Definition 6** (Membership of an operation in a block). Operation  $e$  is a member of block  $b$  in:

- leaf node  $n$ , if  $e$  belongs to  $n.ops[b's\_index]$ .
- internal node  $n$ , if  $e$  is a member of  $s$  that  $s$  is a subblock of  $b$ .

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 **BSearch()**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.

## 7 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.

#### ◇ Structures

##### ► Node

- \**Node* left, right, parent : initialized when creating the tree.
- *BlockList*
- *int* head= 1: #blocks in blocks. blocks[0] is a block with all integer fields equal to zero.
- *int* numpropagated= 0 : # groups of blocks that have been propagated from the node to its parent.

##### ► Block

- *int* group : the value read from numpropagated when appending this 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<sub>enq</sub>, sum<sub>deq</sub> : # enqueue, dequeue operations in the prefix for the 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> : # enqueue operations in the prefix for left.blocks[end<sub>left</sub>]
- *int* sum<sub>deq-left</sub> : # dequeue operations in the prefix for left.blocks[end<sub>left</sub>]
- *int* sum<sub>enq-right</sub> : # enqueue operations in the prefix for right.blocks[end<sub>right</sub>]
- *int* sum<sub>deq-right</sub> : # dequeue operations in the prefix for right.blocks[end<sub>right</sub>]

##### ► 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)
                ▷ Returns the the response to the  $D_{root,b,i}$ .
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
                ▷  $E_e(root)$  is the response.
223:     return root.GetENQ(root.DSEARCH(e, b))
224:   end if
225: end FINDRESPONSE
```

---

---

**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:   <new, npleft, npright>= CREATEBLOCK(h)    ▷ npleft, npright are the
values read from the children's numpropagated field.
312:   if new.num==0 then return true             ▷ The block contains nothing.
313:   else if blocks.tryAppend(new, h) then
314:     for each dir in {left, right} do
315:       CAS(dir.super[npdir], null, h)    ▷ Write would work too.
316:       CAS(dir.numpropagated, npdir, npdir+1)
317:     end for
318:     CAS(head, h, h+1)
319:     return true
320:   else
321:     CAS(head, h, h+1)    ▷ Even if another process wins, help
to increase the head. The winner might have fallen sleep before increasing
head.
322:     return false
323:   end if
324: end REFRESH

↪ Precondition: blocks[start..end] contains a block with field f ≥ i

325: int BSEARCH(field f, int i, int start, int end)
▷ Does binary search for the value
i of the given prefix sum field. Returns the index of the leftmost block in
blocks[start..end] whose field f is ≥ i.

326: end BSEARCH
```

---

---

**Algorithm Root**

---

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

801: <int, int> DSEARCH(int e, int end)    ▷ Returns <b,i> if  $E_e(\text{root}) = E_i(\text{root}, b)$ .
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.BSearch(sumenq, e, start, end)
807:   i= e- root.blocks[b-1].sumenq
808:   return <b,i>
809: end DSEARCH
```

---

---

**Algorithm Node**

---

$\rightsquigarrow$  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.BSEARCH(sumenq, i+blocks[b-1].sumenq-left, blocks[b-1].endleft+1, blocks[b].endleft)
406:     return left.GETENQ(subBlock, i)
407:   else
408:     i= i-blocks[b].numenq-left
409:     subBlock= right.BSEARCH(sumenq, i+right.blocks[b-1].sumenq-right, blocks[b-1].endright+1, blocks[b].endright)
410:     return right.GETENQ(subBlock, i)
411:   end if
412: end GETENQ
```

$\rightsquigarrow$  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.BSEARCH(sumdeq-dir, i+blocks[b-1].sumdeq, super[blocks[b].group]-p, super[blocks[b].group]+p)
▷ superblock's group has at most  $p$  difference with the value stored in  $\text{super}[]$ .
419:     if dir is left then
420:       i+= blocks[b-1].sumenq-blocks[superBlock-1].sumenq-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].sumenq-blocks[superBlock-1].sumenq-right ▷ consider the enqueues in the previous blocks from the right child
424:       i+= blocks[superBlock].numdeq-left ▷ consider the dequeues from the right child
425:     end if
426:     return this.parent.INDEXDEQ(superBlock, i)
427:   end if
428: end INDEXDEQ
```

---

---

**Algorithm Leaf**

---

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

---

---

**Algorithm BlockList**

---

$\triangleright$  : Supports two operations  $\text{blocks.tryAppend}(\text{Block } b)$ ,  $\text{blocks}[i]$ . Initially empty, when  $\text{blocks.tryAppend}(b, n)$  returns true  $b$  is appended to  $\text{blocks}[n]$  and  $\text{blocks}[i]$  returns  $i$ th block in the blocks. If some instance of  $\text{blocks.tryAppend}(b, n)$  returns false there is a concurrent instance of  $\text{blocks.tryAppend}(b', n)$  which has returned true.  $\text{blocks}[0]$  contains an empty block with all fields equal to 0 and  $\text{end}_{\text{left}}$ ,  $\text{end}_{\text{right}}$  pointers to the first block of the corresponding children.

*block[]* blocks: array of blocks  
*int[]* super:  $\text{super}[i]$  stores an approximate index of the superblock of the blocks in blocks whose group field have value  $i$ .

```
701: boolean TRYAPPEND(block blk, int n)
702:   return CAS(blocks[n], null, blk)
703: end TRYAPPEND
```

---

## 8 Proof of Linearizability

**TEST** Fix the logical order of definitions (cyclic references).

**TEST** Is it better to show  $\text{ops}(\text{EST}_n, t)$  with  $\text{EST}_n, t$ ?

**Question** A good notation for *the index of the b*?

**Question** How to remove the notion of time? To say  $\text{pre}(n, i)$  contains  $n.\text{blocks}[0..i]$  instead of  $\text{EST}(n, t)$  which  $\text{head}=i$  at time  $t$ . Is it good? Furthermore, can we remove the notion of established blocks?

**Definition 7** (Block). A block is an object storing some statistics, as described in Algorithm Queue. A block in a node's blocklist implicitly represents a set of operations. If  $n.\text{blocks}[i] = b$  we call  $i$  the *index* of block  $b$ . Block  $b$  is before block  $b'$  in node  $n$  if and only if the index of the  $b$  is smaller than the index of the  $b'$ 's. For a block in a `BlockList` we define *the prefix for the block* to be the blocks in the `BlockList` up to and including the block.

**Lemma 8** (head Increment). *Let  $R$  be an instance of Refresh on node  $n$  that reaches Line 313. After  $R$  terminates  $n.\text{head}$  is greater than  $h$ , the value read in line 310 of  $R$ .*

*Proof.* If Line 318 or 321 are successful then the claim holds, otherwise another process has incremented the head from  $h$  to  $h+1$ .  $\square$

**Invariant 9** (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  $i < h$ .

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

In the algorithm, some value is appended to  $n.\text{blocks}[]$  by writing into  $n.\text{blocks}[\text{head}]$  only in Line 313. Writing into  $n.\text{blocks}[\text{head}]$  preserves the invariant, since the claim does not talk about  $n.\text{blocks}[\text{head}]$ . The value of  $n.\text{head}$  is modified only in lines 318 and 321. Depending on whether the `TryAppend()` in Line 313 succeeded or not, we show that the claim holds after the increment of  $n.\text{head}$  in either case. If  $n.\text{head}$  is incremented to  $h$  it is sufficient to show  $n.\text{blocks}[h] \neq \text{null}$  to prove the invariant still holds. In the first case the process applied a successful `TryAppend(new, h)` in line 314, which means  $n.\text{blocks}[h]$  is not null anymore. Note that whether 318 or 318 return true or false, after they finish we know that  $n.\text{head}$  has been incremented from the value read in Line 310 (Lemma 8). The failure case is also the same since it means some non-null value has been written into  $n.\text{blocks}[\text{head}]$  by some process.  $\square$

*Explain More*

**Lemma 10** (headProgress).  *$n.\text{head}$  is non-decreasing over time. 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.* The first claim follows trivially from the pseudocode since  $n.\text{head}$  is only incremented in the pseudocode in lines 318 and 321 of `Refresh()`.

Consider the block  $b$  written into  $n.\text{blocks}[i]$  by `TryAppend()` at Line 313. It is created by the `CreateBlock(i)` called at Line 311. 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 9. Thus the `CreateBlock(i-1)` that creates  $b'$  terminates before `CreateBlock(i)` that creates  $b$  is invoked. The value written into  $b.\text{end}_{\text{left}}$  at Line 333 of `CreateBlock(i)` was read from  $n.\text{left.head}-1$  at Line 331 of `CreateBlock(i)`. Similarly, the value in  $n.\text{blocks}[i-1].\text{end}_{\text{left}}$  was read from  $n.\text{left.head}-1$  during the call to `CreateBlock(i-1)`. Since  $n.\text{left.head}$  is non-decreasing  $b'.\text{end}_{\text{left}} \leq b.\text{end}_{\text{left}}$ . The proof for  $\text{end}_{\text{right}}$  is similar.  $\square$

**Definition 11** (Subblock). Block  $b$  is a *direct subblock* of  $n.\text{blocks}[i]$  if it is in  $n.\text{left.blocks}[n.\text{blocks}[i-1].\text{end}_{\text{left}}+1..n.\text{blocks}[i].\text{end}_{\text{left}}] \cup n.\text{right.blocks}[n.\text{blocks}[i-1].\text{end}_{\text{right}}+1..n.\text{blocks}[i].\text{end}_{\text{right}}]$ . Block  $b$  is a subblock of  $n.\text{blocks}[i]$  if  $b$  is a direct subblock of  $n.\text{blocks}[i]$  or a subblock of a direct subblock of  $n.\text{blocks}[i]$ .

**Corollary 12** (No Duplicates). *If  $op$  is in  $n.\text{blocks}[i]$  then there is no  $j \neq i$  such that  $op \in \text{ops}(n.\text{blocks}[j])$ .*



*Proof.* Operation `op` is invoked only one time in an execution because every operations invoked is distinct. Since there is node `n` which `op` is in two different blocks of `n`, there is node `n'` that is the lowest height node in the tree that contains `op` in two of its blocks `b1, b2`. By Definition 11, `b1` and `b2` have distinct subblocks(not only direct subblocks) and since `op` is in only one leaf block, then it cannot be in both `b1` and `b2`.  $\square$

**Definition 13** (Superblock). Block `b` is *direct superblock* of block `c` if `c` is a direct subblock of `b`. Block `b` is *superblock* of block `c` if `c` is a subblock of `b`.

**Definition 14** (Operations of a block). A leaf block `b` in a leaf represents `enqueue(x)` if `b.element=x≠null`. Else if `b.element=null` `b` represents a `dequeue()`. The set of operations of block `b` are the operations in the subblocks of `b`. We denote the set of operations of block `b` by `ops(b)`.

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`. We also say `b` contains `op` if `op∈ops(b)`.

**Definition 15.** A block `b` in `n.blocks` is *established* at time `t` if `n.head>` index of `b` at time `t`.  $EST_{n, t}$  is the set of established blocks of node `n` at time `t`.

**Observation 16.** Once a block `b` is written in `n.blocks[i]` then `n.blocks[i]` never changes.

**Lemma 17.** Every block has at most one direct superblock.

*Proof.* To show this we are going to refer to the way `n.blocks[]` is partitioned while propagating blocks up to `n.parent`. `n.CreateBlock(i)` merges the blocks in `n.left.blocks[n.blocks[i-1].endleft..n.blocks[i].endleft]` and `n.right.blocks[n.blocks[i-1].endright..n.blocks[i].endright]` (Lines 331, 332). Since `endleft, endright` are non-decreasing (`n.blocks[i].endleft|right>n.blocks[i-1].endleft|right`), so the range of the subblocks of `n.blocks[i]` which is `(n.blocks[i-1].enddir+1..n.blocks[i].enddir)` does not overlap with the range of the subblocks of `n.blocks[i-1]`.  $\square$

**Lemma 18** (establishedOrder). If time `t < time t'`, then  $ops(EST_{n, t}) \subseteq ops(EST_{n, t'})$ .

*Proof.* Blocks are only appended (not modified) with CAS to `n.blocks[n.head]` and `n.head` is non-decreasing, so the set of operations in established blocks of a node can only grow.  $\square$

*useless?*

► Processes are numbered from 1 to  $p$  and leaves of the tree are assigned from left to right. We will show in Lemma 29 that there is at most one operation from each process in a given block.

**Definition 19** (Ordering of operations inside the nodes). • The prefix of an operation  $op$  in the sequence of operations  $S$  is the sequence of operations strictly before  $op$ .

- $E(n, b)$  is the sequence of enqueue operations in  $\text{ops}(n.\text{blocks}[b])$  defined recursively as follows.  $E(\text{leaf}, b)$  is the single enqueue operation in  $\text{ops}(\text{leaf}.\text{blocks}[b])$  or an empty sequence if  $\text{leaf}.\text{blocks}[b].\text{num}_{\text{enq}}=0$ . If  $n$  is an internal node, then

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

- $E_i(n, b)$  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  $\text{ops}(n.\text{blocks}[b])$  defined recursively as follows.  $D(\text{leaf}, b)$  is the single dequeue operation in  $\text{ops}(\text{leaf}.\text{blocks}[b])$  or an empty sequence if  $\text{leaf}.\text{blocks}[b].\text{num}_{\text{deq}}=0$ . If  $n$  is an internal node, then

$$D(n, b) = D(n.\text{left}, n.\text{blocks}[b-1].\text{end}_{\text{left}} + 1) \cdot D(n.\text{left}, n.\text{blocks}[b-1].\text{end}_{\text{left}} + 2) \cdots D(n.\text{left}, n.\text{blocks}[b].\text{end}_{\text{left}}) \cdot \\ D(n.\text{right}, n.\text{blocks}[b-1].\text{end}_{\text{right}} + 1) \cdot D(n.\text{right}, n.\text{blocks}[b-1].\text{end}_{\text{right}} + 2) \cdots D(n.\text{right}, n.\text{blocks}[b].\text{end}_{\text{right}})$$

- $D_i(n, b)$  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 20** (Linearization).  $L = E(\text{root}, 1).D(\text{root}, 1).E(\text{root}, 2).D(\text{root}, 2).E(\text{root}, 3).D(\text{root}, 3) \cdots$

► In the non-root nodes, we only need ordering of enqueues and dequeues among the operations of their own type. Since `GetENQ()` only searches among enqueues and `IndexDEQ()` works with dequeues.

**Lemma 21** (trueRefresh). *Let  $t_i$  be the time an instance  $R$  of  $n.Refresh()$  is invoked and  $t_t$  be the time it terminates. If the  $TryAppend(new, s)$  of  $R$  returns **true**, then  $ops(EST_{n.left, t_i}) \cup ops(EST_{n.right, t_i}) \subseteq ops(EST_n, t_t)$ .*

*Proof.* Since  $TryAppend$  returns **true** a block **new** is written into  $n.blocks[h]$  in Line 313.

We show  $ops(EST_{n.left, t_i}) \subseteq ops(EST_n, t_t)$ . Let  $h$  be the value  $n.Refresh()$  reads from  $n.head$  at line 310,  $h_{left,i}$  be the value of  $n.left.head$  at  $t_i$  and  $h_{left,read}$  be the value read from  $n.left.head-1$  at line 331.  $end_{left}$  field of the block returned by  $CreateBlock(i)$  is  $h_{left,read}$ . By lines 332 and 331 the **new** block in  $n.blocks[h]$  contains  $n.left.blocks[n.blocks[h-1].end_{left}+1..h_{left,read}]$ . Since  $left.head$  is read after  $t_i$  then  $h_{left,read} > h_{left,i}$  which means  $ops(EST_{n.left, t_i}) \subseteq ops(n.left.blocks[0..h_{left,read}])$ . After the successful  $TryAppend$  in line 313 we know all blocks in  $n.left.blocks[0..h_{left,read}-1]$  are subblocks of  $n.blocks[0..h]$  by the definition of subblock. At  $t_t$  we have  $n.head > h$  by Lemma 10. So  $n.blocks[1..h]$  are in  $EST_{n,t_t}$  by definition of  $EST$ . Note that after line 321 we are sure that the **head** is incremented by Lemma 8) which means  $n.head = h+1$  at  $t_t$  so the new block is established at  $t_t$  and the new block contains the new operations which is what we wanted to show. The proof for  $ops(EST_{n.right, t_i}) \subseteq ops(EST_n, t_t)$  is the same.  $\square$

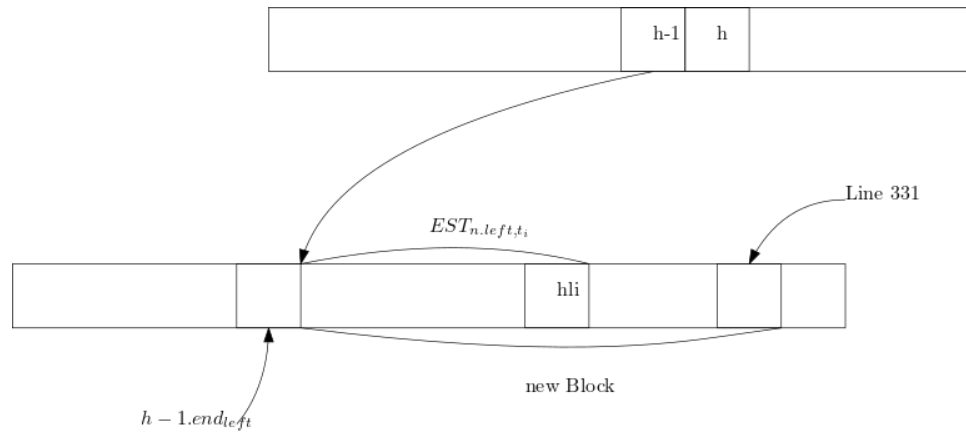


Figure 19: New established operations of the left child are in the new block.

**Lemma 22** (Stronger True Refresh). *Let  $t_i$  be the time an instance of  $n.Refresh()$  read the head (Line 310) and  $t_t$  be the time its  $TryAppend(new, s)$  terminates with and returns **true** (Line 313). We have  $ops(EST_{n.left, t_i}) \cup ops(EST_{n.right, t_i}) \subseteq ops(n.blocks)$ .*

**Definition 23.** An instance of  $Refresh()$  is successful iff its  $TryAppend(new, s)$  terminates with and returns **true**.

**Definition 24.** Let  $R_1 t$  be the time  $R_1$  is invoked and  $t_{R_2}$  be the time  $R_2$  terminates.  $line t$  is the immediate time before running Line  $line$ .  $t_{line}$  is the immediate time after running Line  $line$ .  $line t^{op}$  is the immediate time before running Line  $line$  of operation  $op$ .  $t_{line}^{op}$  is the immediate time after running Line  $line$  of operation  $op$ .

**Lemma 25** (Double Refresh). *Consider two consecutive instances  $R_1, R_2$  of **Refresh()** on internal node  $n$  by a process  $p$ . If  $R_1$  and  $R_2$  both fail and return false, then we have  $ops(EST_{n.left, R_1 t}) \cup ops(EST_{n.right, R_1 t}) \subseteq ops(EST_n, t_{R_2})$ .*

*Proof.*

If  $R_2$  reads some value greater than  $i + 1$  in Line 310 it means a successful instance of **Refresh()** performed its Line 310 after  $t_{310}^{R_1}$  and finished its Line 318 or 321 before  $t_{310}^{R_2}$ , from Lemma 22 by the end of this instance  $ops(EST_{n.left, t_1}) \cup ops(EST_{n.right, t_1})$  has been propagated.

Let  $R_1$  read  $i$  and  $R_2$  read  $i + 1$  from Line 310. As  $R_2$ 's **TryAppend()** returns false, there is another successful instance  $R'_2$  of  $n.Refresh()$  that has done **TryAppend()** successfully into  $n.blocks[i+1]$  before  $R_2$  tries to append. Since  $R'_2$  creates the block after reading the value  $i + 1$  from  $n.head$  (Line 310) and  $R_1$  reads the value  $i$  from  $n.head$  and the  $head$ 's value is increasing by Lemma 10 then  $t_{R'_2 310} > t_{R_1 310} >_{R_1} t$  (see Figure 20). By Lemma 22 after  $R'_2$ 's CAS ( $t_{313}^{R'_2}$ ) we have  $ops(EST_{n.left, t_1}) \cup ops(EST_{n.right, t_1}) \subseteq ops(n.blocks)$ . Also by Lemma 8 on  $R_2$  the value of  $n.head$  head is more than  $i + 1$  after  $R'_2$  terminates, so the block appended by  $R'_2$  to  $n$  is established by then ( $n.head \geq i + 2 > i + 1$ ). To summarize,  $R_1 t$  is before  $R'_2$ 's read of  $n.head$  ( $t_{310}^{R'_2}$ ) and  $R'_2$ 's successful CAS is before  $R_2$ 's termination. So, by Lemma 22,  $ops(EST_{n.left, t_1}) \cup ops(EST_{n.right, t_1}) \subseteq ops(EST_n, t_2)$ .  $\square$

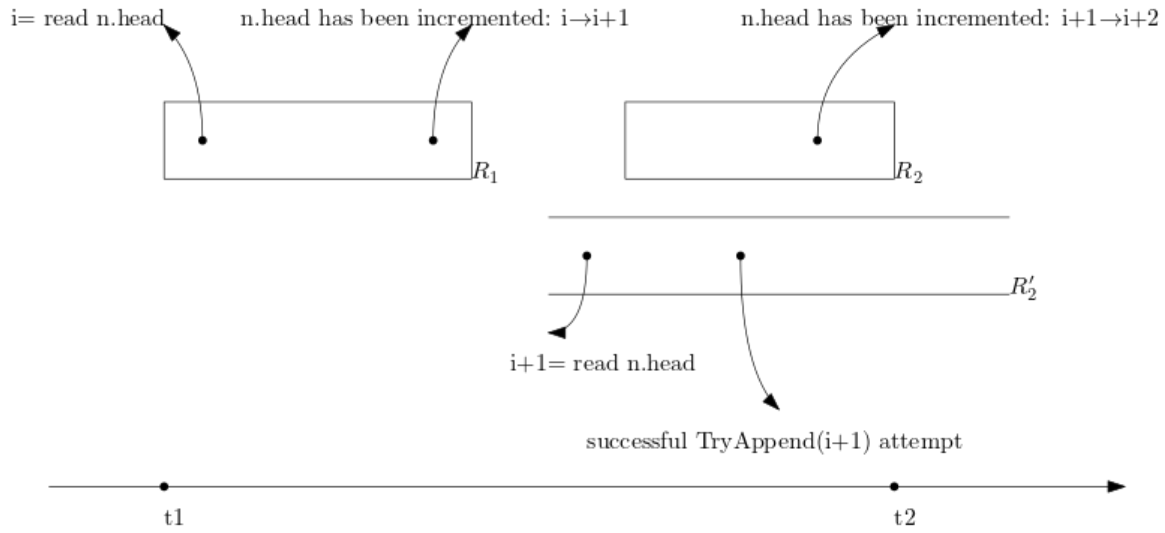


Figure 20:  $R_1 t < t_{310}^{R_1} < \text{incrementing } n.head \text{ from } i \text{ to } i+1 < t_{310}^{R'_2} < t_{313}^{R'_2} < \text{incrementing } n.head \text{ from } i+1 \text{ to } i+2 < t_{R_2}$

**Corollary 26.**  $ops(EST_{n.left, 302t}) \cup ops(EST_{n.right, 302t}) \subseteq ops(EST_n, t_{303})$

*Proof.* If the first **Refresh()** in line 302 returns true then by Lemma 21 the claim holds. Also if first **Refresh()** failed and the second **Refresh()** succeeded the claim still holds by Lemma 21. Finally, if both failed the claim is satisfied by Lemma 25.  $\square$

**Corollary 27** (Propagate Step). *All operations in  $n$ 's children's established blocks before running line 302 of a **Propagate** routine are guaranteed to be in  $n$ 's established blocks after line 303.*

*Proof.* If 302 or 303 succeed, the claim is true by Lemma 21. Otherwise Lines 302 and 303 satisfy the preconditions of Lemma 25.  $\square$

**Corollary 28.** *After **Append**(blk) finishes  $\text{ops}(\text{blk}) \subseteq \text{ops}(\text{root.blocks}[x])$  for exactly one  $x$ .*

*Proof.* After **Append**(blk)'s termination, blk is in **root.blocks** since blk is established in the leaf it has been added to. By applying Lemma 27 inductively it is propagated up to the root. Finally Lemma 12 shows only one block in the root contains blk.  $\square$

**Lemma 29** (Block Size Upper Bound). *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  $op_2$ . A process cannot invoke more than one operations concurrently, so  $op_1$  has to be finished before  $op_2$ . By Corollary 28, before appending  $op_2$  to the tree  $op_1$  exists in every node on the path from  $p$ 's leaf to the root, because  $op_1$ 's **Append** is finished before  $op_2$ 's **Append** starts. So, there is some block  $b'$  before  $b$  in  $n$  containing  $op_1$ . Existence of  $op_1$  in  $b$  and  $b'$  contradicts Lemma 12.  $\square$

**Lemma 30** (Subblocks Upperbound). *Each block has at most  $p$  direct subblocks.*

*Proof.* The claim follows directly from Lemma 29 and the observation that each block appended to the tree contains at least one operation, due to the test on Line 312. We can also see the blocks in the leaves have exactly one operation in the **Enqueue()** and **Dequeue()** routines.  $\square$

**Lemma 31** (Get correctness). *If  $n.blocks[b].num_{enq} \geq i$  then  $n.GetENQ(b, i)$  returns the element enqueued by  $E_i(n, b)$ .*

*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 by the hypothesis, only  $n.GetENQ(b, 1)$  can be called and only when  $n.blocks[b]$  contains an enqueue. At Line 403,  $n.GetENQ(b, 1)$  returns the element of the enqueue operation stored in the  $b$ th block of leaf  $n$ .

For the induction step we prove  $n.GetENQ(b, i)$  returns  $E_i(n, b)$ , assuming  $n.child.GetENQ(subblock, i)$  returns  $E_i(n.child, b)$ . We argue that Line 404 correctly decides whether the  $i$ th enqueue in  $b$ th block of internal node  $n$  is in the left child or right child subblocks of  $n.blocks[b]$ . From Definition 19 of  $E(n, b)$  we know enqueue operations in a block are ordered from left to right and since the leaves of the tree are ordered by process id from left to right, thus operations from the left subblocks come before operations from the right subblocks in a block (See Figure 21). Furthermore the  $num_{enq-left}$  field in  $n.blocks[b]$  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  $i$  is greater than  $b.num_{enq-left}$ . Otherwise we should search for the  $i$ th enqueue in the left child. By definition 11 and 14 we need to search in subblocks of  $n.blocks[b]$  from the range  $n.left.blocks[n.blocks[i-1].end_{left}+1..n.blocks[i].end_{left}] \cup n.right.blocks[n.blocks[i-1].end_{right}+1..n.blocks[i].end_{right}]$ .

If the  $i$ th enqueue of  $n.blocks[b]$  is in the left child it would be  $i$ th enqueue in  $n.left.blocks[n.blocks[i-1].end_{left}+1..n.blocks[i].end_{left}]$  by Definition 11. Also, we know there are  $eb = n.blocks[b-1].sum_{enq-left}$  enqueues in the blocks before this range, so  $E_i(n, b)$  is  $E_{i+eb}(n.left)$  which is  $E_{i'}(n.left, b')$  for some  $b'$  and  $i'$ . We can compute  $b'$  and then search for  $i + eb$ th enqueue in  $n.left$ , where  $i'$  is  $i+eb-n.left.blocks[b'-1].sum_{enq}$ . The parameters in Line 405 are for searching  $E_{i+eb}(n.left)$  in  $n.left.block$  in the expected range of blocks, so this BSearch returns the index of the subblock containing  $E_i(n, b)$ .

Otherwise the enqueue we are looking for is in the right child. Then, there are  $n.blocks[b].num_{enq-left}$  enqueues ahead of it in  $n.blocks[b]$  but not in  $n.right.blocks[n.blocks[i-1].end_{right}+1..n.blocks[i].end_{right}]$ . So we need to search for  $i - n.blocks[b].num_{enq-left} + n.blocks[b-1].sum_{enq-right}$  (Line 409). Other parameters for the left child are chosen similarly to the way they were chosen for the right child.

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

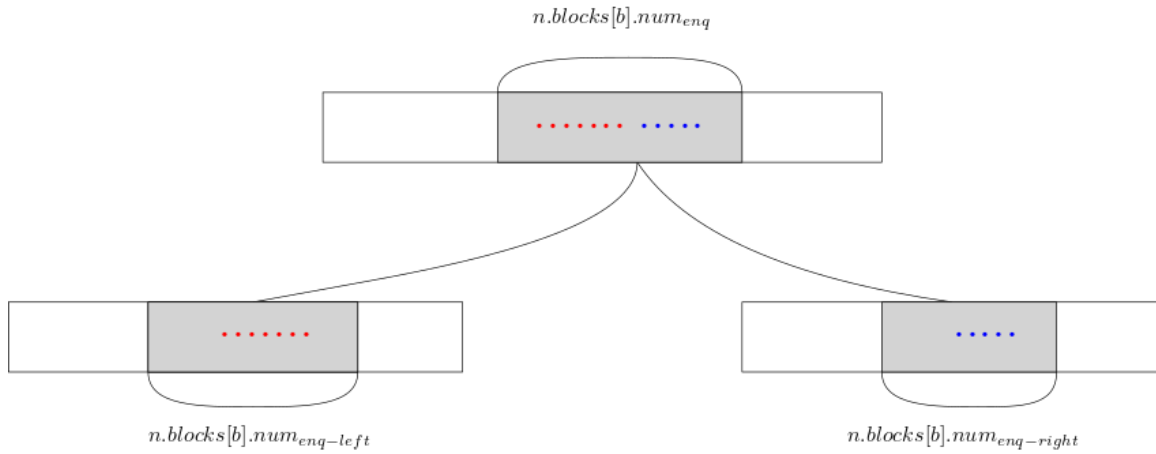


Figure 21: The number and ordering of the enqueue operations propagated from the left and the right child to  $n.blocks[b]$ . Enqueue operations from the left subblocks (colored red), are ordered before the enqueue operations from the right child (colored blue).

**Lemma 32** (DSearch correctness). Assume  $\text{root.blocks}[\text{end}].\text{sum}_{\text{enq}} \geq e$  and  $E_e(\text{root})$ 's element is the response to some `Dequeue()` operation in  $\text{root.blocks}[\text{end}]$ .  $\text{DSearch}(e, \text{end})$  returns  $\langle b, i \rangle$  such that  $E_i(\text{root}, b) = E_e(\text{root})$ .

*Proof.* It is trivial to see that the doubling search from  $\text{root.blocks}[\text{end}]$  to  $\text{root.blocks}[0]$  will find  $E_e(\text{root})$  eventually. Because  $\text{root.blocks}[].\text{sum}_{\text{enq}}$  is an increasing value from 0 to some value greater than  $e$ . So there is a  $b$  that  $\text{root.blocks}[b].\text{sum}_{\text{enq}} > e$  but  $\text{root.blocks}[b-1].\text{sum}_{\text{enq}} < e$ .

First we show  $\text{end} - b \leq 2 \times (\text{root.blocks}[b].\text{size} + \text{root.blocks}[\text{end}].\text{size} + 1)$ . From line 312, we know that size of the every block in the tree is greater than 0. So each block in  $\text{root.blocks}[b..\text{end}]$  contains at least one `Enqueue` or at least one `Dequeue`. Suppose there were more than  $\text{root.blocks}[b].\text{size}$  `Dequeues` in  $\text{root.blocks}[b+1..\text{end}-1]$ . Then the queue would become empty at some point after  $\text{blocks}[b]$ 's last operations and before  $\text{root.blocks}[\text{end}]$ 's first operation. Which means the response to a `Dequeue` in  $\text{root.blocks}[\text{end}]$  could not be in  $E(n, b)$ . Furthermore since the size of the queue would become  $\text{root.blocks}[\text{end}].\text{size}$  after the  $\text{root.blocks}[\text{end}]$ , there cannot be more than  $\text{root.blocks}[b].\text{size} + \text{root.blocks}[\text{end}].\text{size}$  `Enqueues`. Because there can be at most  $\text{root.blocks}[b].\text{size}$  `Dequeues` and the final size is  $\text{root.blocks}[\text{end}].\text{size}$ . Overall there can be at most  $2 \times \text{root.blocks}[b].\text{size} + \text{root.blocks}[\text{end}].\text{size}$  operations in  $\text{root.blocks}[b+1..\text{end}-1]$  and since each block size is  $\geq 1$  thus 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  $\text{end} - b \leq 2 \times \text{root.blocks}[b].\text{size} + \text{root.blocks}[\text{end}].\text{size} + 1$ . See Figure ??.

Now that we know there are at most  $\text{root.blocks}[b].\text{size} + \text{root.blocks}[\text{end}].\text{size}$  blocks in between  $\text{root.blocks}[b]$  and  $\text{root.blocks}[\text{end}]$  then with doubling search in  $\Theta(\log(\text{root.blocks}[b].\text{size} + \text{root.blocks}[\text{end}].\text{size}))$  steps we reach  $\text{start} = c$  that the  $\text{root.blocks}[c].\text{sum}_{\text{enq}}$  is less than  $e$  and  $\text{end} - c$  is not more than  $2 \times \text{root.blocks}[b].\text{size} + \text{root.blocks}[\text{end}].\text{size}$ . Beause otherwise, then  $(\text{end} - c)/2$  satisfied the  $\text{root.blocks}[(\text{end} - c)/2].\text{sum}_{\text{enq}} < e$ . In line 804 the differenece between  $\text{end}$  and  $\text{start}$  is doubled. See Figure 22.

After computing  $b$ , the value  $i$  is computed via the definition of  $\text{sum}_{\text{enq}}$  in constant time (Line 807). So the routine non constant part is the binary search which takes  $\Theta(\log \text{root.blocks}[b].\text{size} + \text{root.blocks}[\text{end}].\text{size})$  steps from the first paragraph.

□

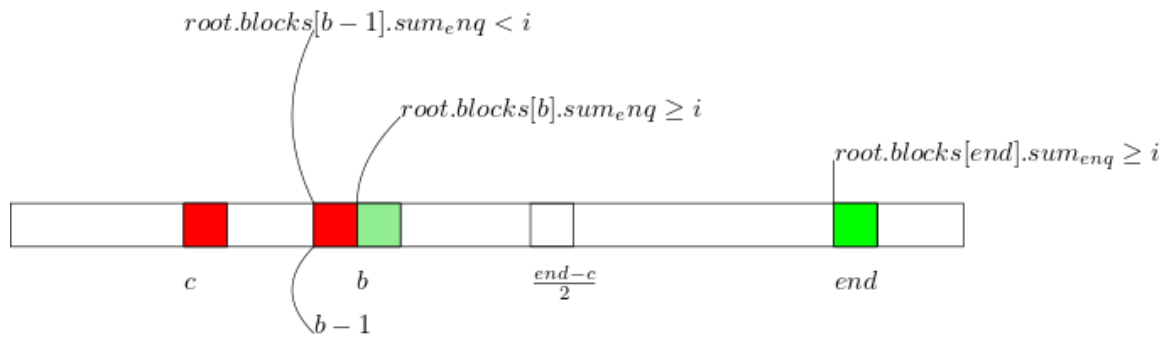


Figure 22: Distance relations between  $b, c, \text{end}$

**Lemma 33.** Let  $n.propagates$  be the number of groups of blocks that have been propagated from node  $n$  to its parent (successful  $n.parent.Refresh()$ ). We have  $num_{propagated} \leq n.propagates \leq num_{propagated} + p$ .  $p$  is the number of processes.

*Proof.*  $num_{propagated}$  is incremented after propagating (Line 316). Since maybe some process falls sleep before incrementing  $num_{propagated}$  it may be behind by  $p$ . □

**Lemma 34.**  $super[]$  preserves order from child to parent; i.e. if in node  $n$  block  $b$  is before  $c$  then  $b.group \leq c.group$

*Proof.* Line 329. Since  $num_{propagated}$  is increasing. □

**Lemma 35.** Let  $b, c$  be in node  $n$ , if  $b.group \leq c.group$  then  $super[b.group] \leq super[c.group]$

*Proof.* Line 315. □

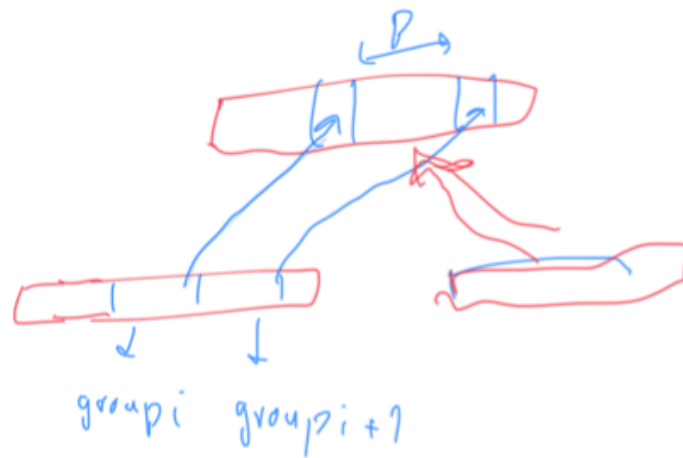
**Lemma 36.** The number of the blocks with  $group=i$  in a node is  $\leq p$ .

*Proof.* For the sake of simplicity we assumed all the blocks are propagated from the left child. □



**Lemma 37.**  $super[i+1] - super[i] \leq p$

*Proof.* In a Refresh with successful CAS in line 46,  $super$  and  $counter$  are set for each child in lines 48,49. Assume the current value of the counter in node  $n$  is  $i+1$  and still  $super[i+1]$  is not set. If an instance of successful  $Refresh(n)$  finishes  $super[i+1]$  is set a new value and a block is added after  $n.parent[super[i]]$ . There could be at most  $p$  successful unfinished concurrent instances of  $Refresh()$  that have not reached line 49. So the distance between  $super[i+1]$  and  $super[i]$  is less than  $p$ . □



**Lemma 38** (super property). If  $super[i] \neq null$  in node  $n$ , then  $super[i]$  is the index of the superblock of a block with  $time=i$  in  $n.parent.blocks$ .

**Lemma 39.** Superblock of  $b$  is within range  $\pm 2p$  of the  $super[b.group]$ .



*Proof.*  $\text{super}[i]$  is the index of the superblock of a block containing block  $b$ , followed by Lemma 38.  $\text{super}(b)$  is the real superblock of  $b$ .  $\text{super}(t)$  is the index of the superblock of the last block with time  $t$ . If  $b.\text{time}$  is  $t$  we have:

$$\text{super}[t] - p \leq \text{super}[t - 1] \leq \text{super}(t - 1) \leq \text{super}(b) \leq \text{super}(t + 1) \leq \text{super}(t + 1) \leq \text{super}[t] + p$$

□

**Lemma 40.** *Search in each level of  $\text{IndexDeq}()$  takes  $O(\log p)$  steps.*

*Proof.* Show preconditions are satisfied and the range is  $p$ .

□

**Lemma 41** (Computing SuperBlock). *For the  $\text{superblock}$  value computed in line 418 of  $\text{n.IndexDEQ}(b, i)$  we have  $\text{n.parent.blocks}[\text{superblock}]$  contains  $D_{n,b,i}$ .*

*Proof.* First we show the value read for  $\text{super}[b.\text{group}]$  in line 418 is not null. Values  $\text{np}_{\text{dir}}$  read in lines 337,  $\text{super}$  are set before incrementing in lines 315,316. So before incrementing  $\text{num}_{\text{propagated}}$ ,  $\text{super}[\text{num}_{\text{propagated}}]$  is set so it cannot be null while reading. Then by Lemma 39 if we search in the range  $p$ , we can find the superblock.

□

**Lemma 42** (Index correctness). *If  $\text{n.blocks}[b].\text{num}_{\text{deq}} \geq i$  then  $\text{n.IndexDEQ}(b, i)$  returns the rank in  $D(\text{root})$  of  $D_{n,b,i}$ .*

*Proof.* We will prove this by induction on the distance of  $n$  from the  $\text{root}$ . We can see the base case where  $n$  is root is trivial (Line 415). In the non-root nodes  $\text{n.IndexDEQ}(b, i)$  computes the superblock of the  $i$ th Dequeue in the  $b$ th block of  $n$  in  $\text{n.parent}$  by Lemma 41 (Line 418). After that the order in  $D(n.\text{parent}, \text{superblock})$  is computed. Note that by Lemma 29 in each block there is at most one operation from each process and operations of one type are ordered based on the order in the subblocks (See Figure 23). Finally  $\text{index}()$  is called on  $\text{n.parent}$  recursively and it returns the correct response from induction hypothesis. If the operation was propagated from the right child the number of dequeues from the left child are added to it (Line ??), because the left child operations come before the right child operations (Definition 19).

□

*Make sure to show preconditions of all invocation of  $\text{BSearch}$  are satisfied.*

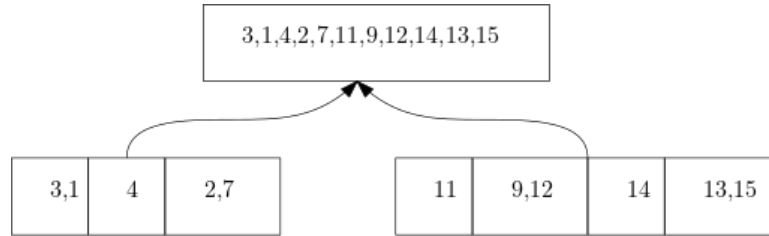


Figure 23: Relation of ordering of operations of a block from its subblocks

**Definition 43.** Assume the operations in  $L$  are applied on an empty queue. If element of `enqueue e` is the response to `dequeue d` then we say  $R(d)=e$ . If  $d$ 's response is `null` (queue is empty) then  $R(d)=\text{null}$ .

**Definition 44.** In an execution on a queue, the dequeue operations that return some value are called *non-null dequeues*.

**Observation 45.** In a sequential execution on a queue,  $k$ th non-null dequeue returns the `element` of  $k$ th enqueue.

**Lemma 46.** `root.blocks[b].size` is the size of the queue if the operations in the prefix for the  $b$ th block in the root are applied with the order of  $L$ .

*Proof.* need to say? :: If the size of a queue is greater than 0 then a `Dequeue()` would decrease the size of the queue, otherwise the size of the queue remains 0. By definition 19 enqueue operations come before dequeue operations in a block in  $L$ .

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`. For  $b=i$  we are going to use the hypothesis for  $b=i-1$ . 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 we can compute the size of the queue after  $b$ th block using with this equality `root.blocks[b].size= root.blocks[b-1].size+ root.blocks[b].sum_enq- root.blocks[b].sum_deq` (Line 342). See Table 4 for an example of running some blocks of operations on an empty queue.  $\square$

**Lemma 47** (Duality of #non-null dequeues and `block.size`). If the operations are applied with the order of  $L$ , the number of non-null dequeues in the prefix for a block  $b$  is `b.sum_enq-b.size`

*Proof.* There are `b.sum_enq` enqueue operations in the prefix for  $b$ , then the size of the queue after the prefix for  $b$  is `#enqs - #non-null dequeues` in the prefix for  $b$ , by Observation 35. So `#non-null dequeues` is `b.sum_enq-b.size`. The correctness of the `block.size` field is shown in Lemma 46.  $\square$

**Lemma 48.**  $R(D_{\text{root},b,i})$  is null iff `root.blocks[b-1].size + root.blocks[b].num_enq- i < 0`.

**Lemma 49** (Computing Response). `FindResponse(b,i)` returns  $R(D_{\text{root},b,i}).\text{element}$ .

*Proof.* First note that by Definition 19 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.

$D_{\text{root},b,i}$  is  $D_{\text{root},\text{root.blocks}[b-1].\text{sum\_deq}+i}$  from the definition 19 and  $\text{sum\_enq}$ .  $D_{\text{root},b,i}$  returns null if `root.blocks[b-1].size + root.blocks[b].num_enq- i < 0` by Lemma 48 (Line 220). Otherwise if it is  $d'$ th non-null dequeue in  $L$  it returns  $d'$ th enqueue by Observation 45. By Lemma 47 there are `root.blocks[b-1].sum_enq - root.blocks[b-1].size` non-null dequeue operations before prefix for `root.blocks[b-1]`. Note that the dequeues in `root.blocks[b]` before the  $i$ th dequeue are non-null dequeues. So the response is  $E_{i-\text{root.blocks}[b-1].\text{size}+\text{root.blocks}[b-1].\text{sum\_deq}}(\text{root})$  (Line 222). See figure 24.

After computing  $e$  we can find  $b,i$  such that  $E_i(\text{root},b) = E_e(\text{root})$  using `DSearch` and then find its `element` using `GetEnq` (Line 223).  $\square$

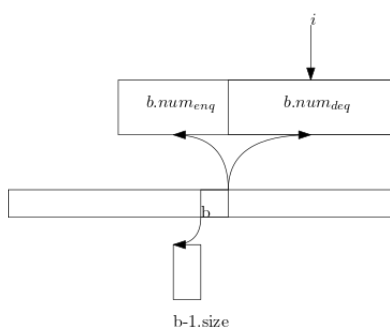


Figure 24: The position of  $E_i(\text{root},b)$ .

	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 4: An example of root blocks fields. Blocks are from left to right and operations in the blocks are also from the left to right.

**Theorem 50** (Main). *The queue implementation is linearizable.*

*Proof.* We choose  $L$  in Definition 19 to be linearization ordering of operations and prove if we linearize operations as  $L$  the queue works consistently.  $\square$

**Lemma 51** (satisfiability).  *$L$  can be a linearization ordering.*

*Proof.* To show this we need to say if in an execution,  $op_1$  terminates before  $op_2$  starts then  $op_1$  is linearized before  $op_2$ . If  $op_1$  terminates before  $op_2$  starts it means  $op_1.\text{Append}()$  is terminated before  $op_2.\text{Append}()$  starts. From Lemma 12  $op_1$  is in `root.blocks` before  $op_2$  propagates so  $op_1$  is linearized before  $op_2$  by Definition 19.

Once some operations are aggregated in one block they will be propagated together up to the root and we can linearize them in any order among themselves. Furthermore in  $L$  we arbitrary choose the order to be by process id, since it makes computations in the blocks faster.  $\square$

**Lemma 52** (correctness). *If operations are applied as  $L$  on a sequential queue, the sequence of the responses would be the same as our algorithm.*

*Proof. Old parts to review* We show that the ordering  $L$  stored in the root, satisfies the properties of a linearizable ordering.

1. If  $op_1$  ends before  $op_2$  begins in  $E$ , then  $op_1$  comes before  $op_2$  in  $T$ .
  - This is followed by Lemma 12. The time  $op_1$  ends it is in root, before  $op_2$ , by Definition 19  $op_1$  is before  $op_2$ .
2. Responses to operations in  $E$  are same as they would be if done sequentially in order of  $L$ .
  - Enqueue operations do not have any response so it does no matter how they are ordered. It remains to prove Dequeue  $d$  returns the correct response according to the linearization order. By Lemma 49 it is deduced that the head of the queue at time of the linearization of  $d$  is computed properly. If the Queue is not empty by Lemma 31 we know that the returning response is the computed index element.

$\square$

**Lemma 53** (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 39) 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 32, 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 30), 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. □

**Lemma 54** (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 313 and one in Lines 318 or 321. □