# 1 Pseudocode

---

**Algorithm** Fields description

---

**Structure**

◇ *Shared*

- *Tree* tree : A binary tree of Nodes. root is a pointer to the root node.

◇ *Local*

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

◇ *Structures*

▶ *Node*

- *Node* left, right, parent : initialized when creating the tree.

- *BlockList* blocks implemented with an array.

- *int* size= 1: #blocks in blocks.

- *int* $\text{num}_{\text{propagated}}$= 0 : # groups of blocks that have been propagated from the node to its parent. Since it is incremented after propagating, it may be behind by 1.

- *int[]* super: super[i] stores an approximate index of the superblock of the blocks in blocks whose group field have value i.

▶ *Leaf* extends *NonRootNode*

- *int* $\text{last}_{\text{done}}$

  Stores the index of the block in the root such that the process that owns this leaf has most recently finished the. A block is finished if all of its operations are finished. enqueue(e) is finished if e is returned by some dequeue() and dequeue() is finished when it computes its response. *put the definitions before the pseudocode*

▶ *Block*      ▷ 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. *put the definitions before the pseudocode*

- *int* group : the value read from $\text{num}_{\text{propagated}}$ when appending this block to the node.

▶ *LeafBlock* extends *Block*

- *Object* element : Each block in a leaf represents a single operation. For enqueue operations element is the input of the enqueue and for dequeue operations it is null.

- *Object* response : stores the response of the operation in the LeafBlock.

- *int* $\text{sum}_{\text{enq}}$, $\text{sum}_{\text{deq}}$ : # enqueue, dequeue operations in the prefix for the block

▶ *InternalBlock* extends *Block*

- *int* $\text{end}_{\text{left}}$, $\text{end}_{\text{right}}$ : index of the last subblock of the block in the left and right child

- *int* $\text{sum}_{\text{enq-left}}$ : # enqueue operations in the prefix for left.blocks[$\text{end}_{\text{left}}$]

- *int* $\text{sum}_{\text{deq-left}}$ : # dequeue operations in the prefix for left.blocks[$\text{end}_{\text{left}}$]

- *int* $\text{sum}_{\text{enq-right}}$ : # enqueue operations in the prefix for right.blocks[$\text{end}_{\text{right}}$]

- *int* $\text{sum}_{\text{deq-right}}$ : # dequeue operations in the prefix for right.blocks[$\text{end}_{\text{right}}$]

▶ *RootBlock* extends *InternalBlock*

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

- *counter* $\text{num}_{\text{finished}}$ : number of finished operations in the block

- *int* order : the index of the block in the BlockList containing the block.

---

*Variable naming:*

- $\text{b}_{\text{op}}$: index of the block containing operation op

- $\text{r}_{\text{op}}$: rank of operation op i.e. the ordering among the operations of its type according to linearization ordering

*Abbreviations:*

- blocks[b].$\text{sum}_{\text{x}}$=blocks[b].$\text{sum}_{\text{x-left}}$+blocks[b].$\text{sum}_{\text{x-right}}$ (for b≥0 and x ∈ {enq, deq})

- blocks[b].sum=blocks[b].$\text{sum}_{\text{enq}}$+blocks[b].$\text{sum}_{\text{deq}}$ (for b≥0)

- blocks[b].$\text{num}_{\text{x}}$=blocks[b].$\text{sum}_{\text{x}}$-blocks[b-1].$\text{sum}_{\text{x}}$

  (for b>0 and x ∈ {∅, enq, deq, enq-left, enq-right, deq-left, deq-right}, blocks[0].$\text{num}_{\text{x}}$=0)

**Algorithm** *Queue*

---

201: *void* Enqueue(*Object* e)  ▷ Creates a block with element e and appends

  it to the tree.

202:   block newBlock= NEW(*LeafBlock*)

203:   newBlock.element= e

204:   newBlock.sum$_{enq}$= leaf.blocks[leaf.size].sum$_{enq}$+1

205:   newBlock.sum$_{deq}$= leaf.blocks[leaf.size].sum$_{deq}$

206:   leaf.Append(newBlock)

207: **end** Enqueue


208: *Object* Dequeue()

209:   block newBlock= NEW(*LeafBlock*)   ▷ Creates a block

  with null value element, appends it to the tree, computes its order among

  operations, then computes and returns its response.

210:   newBlock.element= null

211:   newBlock.sum$_{enq}$= leaf.blocks[leaf.size].sum$_{enq}$

212:   newBlock.sum$_{deq}$= leaf.blocks[leaf.size].sum$_{deq}$+1

213:   leaf.Append(newBlock)

214:   **return** leaf.HelpDequeue()

215: **end** Dequeue

216: *int, int* FindResponse(int i, int b)   ▷ Computes the rank and

  index of the block in the root of the enqueue that is the response of the ith

  dequeue in the root's bth block. Returns <-1,--> if the queue is empty.

217:   **if** root.blocks[b-1].length + root.blocks[b].num$_{enq}$ - i $<$ 0 **then**

218:     **return** <-1,-->

219:   **else**

                              ▷ We call the dequeues that

  return a value *non-null dequeues*. $r$th non-null dequeue returns the element

  of th $r$th enqueue. We can compute # non-null dequeues in the prefix for

  a block this way: #non-null dequeues= length - #enqueues. Note that the

  $i$th dequeue in the given block is not a non-null dequeue.

220:     r$_{enq}$= root.blocks[b-1].sum$_{enq}$- root.blocks[b-1].length + i

221:     **return** <root.blocks.get(enq, r$_{enq}$).order, r$_{enq}$>

222:   **end if**

223: **end** FindResponse

---

**Algorithm** *Node*

301: *void* PROPAGATE()

302:    **if not** REFRESH() **then**

303:        REFRESH()                          ▷ Lemma Double Refresh

304:    **end if**

305:    **if** this **is not** root **then**

306:        parent.PROPAGATE()

307:    **end if**

308: **end** PROPAGATE

309: *boolean* REFRESH()                                        $\boxed{\text{lastLine}}$

310:    s= size                                                 $\boxed{\text{prevLine}}$

311:    <new, $np_{left}$, $np_{right}$>= CREATEBLOCK(h)        ▷ $np_{left}$, $np_{right}$ are the
   values read from the children's $num_{propagated}$ feild.

312:    **if** new.num==0 **then return** true         ▷ The block contains nothing.

313:    **else if** blocks.tryAppend(new, s) **then**

$\boxed{\text{okcas}}$314:        **for each** dir **in** {left, right} **do**

315:            CAS(dir.super[$np_{dir}$], null, h+1)    ▷ Write would work too.

316:            CAS(dir.num$_{propagated}$, $np_{dir}$, $np_{dir}$+1)

317:        **end for**

318:        CAS(size, h, h+1)

319:        **return** true

320:    **else**

321:        CAS(size, h, h+1)                          ▷ Even if another
   process wins, help to increase the size. The winner might have fallen sleep
   before increasing size.

322:        **return** false

323:    **end if**

324: **end** REFRESH

   ↝ Precondition: blocks[start..end] contains a block with field f $\geq$ 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 $\geq$ i.

326: **end** BSEARCH

327: <*Block, int, int*> CREATEBLOCK(*int* i)
                          ▷ Creates a block to be inserted into as $i$th block in
   blocks. Returns the created block as well as values read from each child's
   $num_{propagated}$ field. These values are used for incrementing the children's
   $num_{propagated}$ field if the block was appended to blocks successfully.

328:    block newBlock= NEW(*block*)

329:    newBlock.group= num$_{propagated}$

330:    newBlock.order= i

331:    **for each** dir **in** {left, right} **do**

332:        index$_{last}$= dir.size

333:        index$_{prev}$= blocks[i-1].end$_{dir}$

334:        newBlock.end$_{dir}$= index$_{last}$

335:        block$_{last}$= dir.blocks[index$_{last}$]

336:        block$_{prev}$= dir.blocks[index$_{prev}$]

337:            ▷ newBlock includes dir.blocks[index$_{prev}$+1..index$_{last}$].

338:        this$_{dir}$= dir.num$_{propagated}$

339:        newBlock.sum$_{enq-dir}$= blocks[i-1].sum$_{enq-dir}$ + block$_{last}$.sum$_{enq}$
   - block$_{prev}$.sum$_{enq}$

340:        newBlock.sum$_{deq-dir}$= blocks[i-1].sum$_{deq-dir}$ + block$_{last}$.sum$_{deq}$
   - block$_{prev}$.sum$_{deq}$

341:    **end for**

342:    **if** this **is** root **then**

343:        newBlock.length= max(root.blocks[i-1].length + b.num$_{enq}$ -
   b.num$_{deq}$, 0)

344:    **end if**

345:    **return** <b, $np_{left}$, $np_{right}$>

346: **end** CREATEBLOCK

**Algorithm** Node

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

401: *element* GETENQ(*int* b, *int* i)

402:     **if** this **is** leaf **then**

403:         **return** blocks[b].element

404:     **else if** $i \leq \text{blocks[b].num}_{\text{enq-left}}$ **then**                         ▷ i exists in the left child of this node

405:         subBlock= left.BSEARCH($\text{sum}_{\text{enq}}$, i, $\text{blocks[b-1].end}_{\text{left}}$+1, $\text{blocks[b].end}_{\text{left}}$)      ▷ Search range of left child's subblocks of blocks[b].

406:         **return** left.GET($\text{i-left.blocks[subBlock-1].sum}_{\text{enq}}$, subBlock)

407:     **else**

408:         $i = \text{i-blocks[b].num}_{\text{enq-left}}$

409:         subBlock= right.BSEARCH($\text{sum}_{\text{enq}}$, i, $\text{blocks[b-1].end}_{\text{right}}$+1, $\text{blocks[b].end}_{\text{right}}$)     ▷ Search range of right child's subblocks of blocks[b].

410:         **return** right.GET($\text{i-right.blocks[subBlock-1].sum}_{\text{enq}}$, subBlock)

411:     **end if**

412: **end** GETENQ


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

413: <*int, int*> INDEXDEQ(*int* b, *int* i)                   ▷ Returns the rank of $i$th dequeue in the $b$th block of the node, among the dequeues in the root.

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

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

416:     **else**

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

418:         superBlock= parent.BSEARCH($\text{sum}_{\text{deq-dir}}$, i, super[blocks[b].group]-p, super[blocks[b].group]+p)

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

419:         **if** dir **is** right **then**

420:             $\text{i+= blocks[superBlock].sum}_{\text{deq-left}}$                 ▷ consider the dequeues from the right child

421:         **end if**

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

423:     **end if**

424: **end** INDEX

## Algorithm Leaf

501: *void* APPEND(*block* blk)                                                                    ▷ Append is only called by the owner of the leaf.

502:     size+=1

503:     blk.group= size

504:     blocks[size]= blk

505:     parent.PROPAGATE()

506: **end** APPEND


507: *Object* HELPDEQUEUE()

508:     <$b_{deq}$, $r_{deq}$>= INDEXDEQ(leaf.size, 1)                    ▷ r is the rank among the dequeues of the dequeue of the $b_{deq}$th block in the root containing.

509:     <$r_{enq}$, $r_{enq}$>= FINDRESPONSE($r_{deq}$, $b_{deq}$)  ▷ $r_{enq}$ is the rank of the enqueue whose element is the response to the dequeue in the block containing it and

     $b_{deq}$ is the index of that block of it in the blocklist. If the response is null then $r_{deq}$ is -1.

510:     **if** $r_{enq}$==-1 **then**

511:         output= null

512:         root.blocks[$b_{deq}$].num$_{finished}$.inc()                                                       ▷ shared counter

513:         **if** root.blocks[$b_{deq}$].num$_{finished}$==root.blocks[$b_{deq}$].num **then**

514:             last$_{done}$= $b_{deq}$

515:         **end if**

516:     **else**

517:         output= GETENQ($b_{enq}$, $r_{enq}$)                                                       ▷ getting the reponse's **element**.

518:         root.blocks[$b_{enq}$].num$_{finished}$.inc()

519:         root.blocks[$b_{enq}$].num$_{finished}$.inc()

520:         **if** root.blocks[$b_{deq}$].num$_{finished}$==root.blocks[$b_{deq}$].num **then**

521:             last$_{done}$= $b_{deq}$

522:         **else if** root.blocks[$b_{enq}$].num$_{finished}$==root.blocks[$b_{enq}$].num **then**

523:             last$_{done}$= $b_{enq}$

524:         **end if**

525:     **end if**

526:     **return** output

527: **end** DEQUEUE

---

**Algorithm** BlockList

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

◇ *PBRT implementation*

A persistant red-black tree supporting `append(b, key)`,`get(key=i)`,`split(j)`. `append(b, key)` returns `true` in case successful. Since order, $sum_{enq}$are both strictly increasing we can use one of them for another.

root: pointer to the root of the PBRT

601: *boolean* TRYAPPEND(*block* blk, *int* n)                                                              ▷ adds block b to the `root.blocks[n]`

602:     **if** `root.size%`$p^2$`==0` **then**                                                 ▷ Help every often $p^2$ operations appended to the root.

603:         `Help()`

604:         `CollectGarbage()`

605:     **end if**

606:     `blk.num`$_{finished}$`= 0`

607:     `*oldRoot= &root.blocks.root`

608:     `*newRoot= root.blocks.Append(blk).root`

609:     **return** `CAS(root, oldRoot, newRoot)`

610: **end** TRYAPPEND

◇ *Array implementation*

blocks[]: array of blocks

611: *boolean* TRYAPPEND(*block* blk, *int* n)

612:     **return** `CAS(blocks[n], null, blk)`

613: **end** TRYAPPEND

---

801: *void* HELP                                     ▷ Helps pending operations

802:     **for** `leaf l` **in** `leaves` **do**

803:         `last= l.size-1`   ▷ `l.blocks[last]` can not be `null` because size increases after appending, see lines 503-502.

804:         **if** `l.blocks[last].element==null` **then**   ▷ operation is dequeue

805:             `l.blocks[last].response= l.`HELPDEQUEUE`()`

806:         **end if**

807:     **end for**

808: **end** HELP

809: *void* COLLECTGARBAGE              ▷ Collects the root blocks that are done.

810:     `s=FindMostRecentDone(Root.Blocks.root)`   ▷ Lemma: If block b is done after helping then all blocks before b are done as well.

811:     `t1,t2= RBT.split(order, s)`

812:     `RBTRoot.CAS(t2.root)`

813: **end** COLLECTGARBAGE

814: *Block* FINDMOSTRECENTDONE(b)

815:     **for** `leaf l` **in** `leaves` **do**

816:         `max= Max(l.maxOld, max)`

817:     **end for**

818:     **return** `max`                                       ▷ This snapshot suffies.

819: **end** FINDMOSTRECENTDONE

820: *response* FALLBACK(op i) ▷ *how to use as exception handling? by adding try catch in all the methods reading the root?*

821:     **if** `root.blocks.get(num`$_{enq}$`), i is null` **then**   ▷ this enqueue was already finished

822:         **return** `this.leaf.response(block.order)`

823:     **end if**

824: **end** FALLBACK

---

# 2 Proof of Linearizability

**Definition 1.** If `n.blocks[i]==b` we call `i` the *index* of block `b` in node `n`. Block `b` is before block `b'` in node `n` if and only if `b`'s index is smaller than `b'`'s. Block `b` is propagated to node `n` or set `S` if `b` is in `n.blocks` or `S` or is a subblock of a block in `n.blocks` or `S`.

**Definition 2.** Block `b` in node `n` is in the set *Established(n, t)* if `n.size` is greater than `b`'s index at time $t$.

**Lemma 3** (sizeProgress). `n.size` *is non-decreasing over time.*

**Lemma 4** (headPosition). *The value read in Line 333(`h=n.size`) might be 1 bit behind the first empty block in the node.*

**Lemma 5** (establishedOrder). *If time $t <$ time $t'$, then $Established(n,t) \subseteq Established(n,t')$.*

**Lemma 6** (createBlock). *Suppose `n.CreateBlock( h, x)` is invoked at time $t$. The blocks propagated to $Established(n.left, t)$ and $Established(n.right, t)$ that are not propagated to $Established(n, t)$, are subblock of the block returned by `CreateBlock(n, h, x)`.*

**Lemma 7** (trueRefresh). *Suppose `Refresh(n)`'s `CAS(n.blocks[h], null, new)` returns `true`. Let $t$ be the time `Refresh(n)` is invoked, blocks propgatated to $Established(n.left, t)$ and $Established(n.right, t)$ are propagated to in $Established(n, t)$ after `CAS(n.blocks[h], null, new)`.*

**Lemma 8** (falseRefresh). *If instance $r$ of `Refresh(n)` returns `false`, then there is another successful instance $r'$ of `Refresh(n)` that has performed a successful `CAS(n.blocks[h], null, new)`(Line 49) after Line 43(`h= n.head`) of $r$.*

**Lemma 9** (Double Refresh). *Consider two consecutive instances $r_1, r_2$ of `Refresh(n)` by the same process (Lines 35,36). Let be the time before $r_1$ invoked. After $r_2$'s `CAS` all the blocks propagated to $Established(n.left, t)$ and $Established(n.right, t)$ are in $Established(n, t)$.*
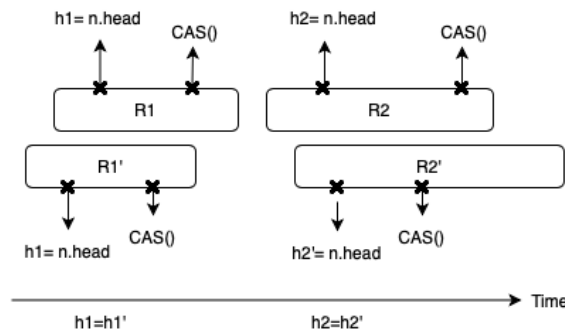
Figure 1: $R_2'$'s CAS is executed after `h1=n.head`.

Block `new` is created of new established subblocks of children of n(Lemma 6, Line 46). If CAS in Line 48 succeeds then by Lemma 7 new established blocks will be in `n`.

**Lemma 10** (Double Refresh). *All operations in `n`'s children's blocks before line 35 are guaranteed to be in `n`'s blocks after Line 37.*

`CreateBlock()` reads blocks in the children that do 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 `Refresh`es 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 35. 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.

**Lemma 11** (Append). *When `Append(op)` is finished, `op` appears exactly once in a block of `root.blocks`.*

**Lemma 12** (Block Size Upper Bound). *Each block in a node contains at most one operation from each processs.*

**Lemma 13** (Subblocks Upperbound). *Each block in a node has at most p subblocks.*

**Definition 14** (Ordering of operations inside a node)**.** ► Note that from Lemma $\overset{\text{blockSize}}{\overline{12 \text{ we know}}}$ there is at most one operation from each process in a given block.

- $E(n, i)$ is the sequence of enqueue operations that are member of `n.blocks[i]` ordered by process id.

- $D(n, i)$ is the sequence of dequeue operations that are member of `n.blocks[i]` ordered by process id.

- $D(n) = D(n, 1).D(n, 2).D(n, 3)...$

- $L = E(root, 1).D(root, 1).E(root, 2).D(root, 2).E(root, 3).D(root, 3)...$

**Theorem 15.** *The queue implementation is linearizable.*

**Lemma 16** (Get)**.** *`Get(n,b,i)` returns ith Enqueue in $E(n, b)$.*

**Lemma 17** (Index)**.** *Index(n,b,i) returns the rank in the $D(root)$ of ith Dequeue in $D(n, b)$.*

**Lemma 18** (Computing SuperBlock)**.** *If `Index(n,b,i)` performs line 101, then `superblock` contains ith Dequeue in bth block of node n.*

**Lemma 19** (Computing Queue's Head)**.** *Let $Q$ be state of the queue if the operations before ith Dequeue in $L(root)$ are applied on the Queue sequentially and $X$ be the head of $Q$. If $Q$ is empty `ComputeHead(i,b)` returns -1, otherwise returns index in $E(root, b)$ of X.*

**Lemma 20** (Validity of `head`)**.** *No two blocks are written in the same index in `n.blocks`.*

**Lemma 21** (Validity of `super` and `counter`)**.** *If `super[i]` $\neq$ `null`, then `super[i]` in node `n` is the superblock of a block with `time=i`.*

**Lemma 22** (Search Ranges)**.** *Preconditions of all invocation of `BSearch` are satisfied.*