

1 Elimination Stack

For simplicity, we assume a value \perp , which is different from NULL and any other value the stack can store. Since NULL is used as a legit return value, representing the value of POP operation (when exchanging values using the elimination array), NULL can not be used to represent an initialization value, different then any stack value. The same holds for a Node, since a NULL node represent an empty stack, the value \perp is used to distinguish between initialization value and empty stack.

For simplicity, we split the RECOVER routine into sub-routines, based on which operation (PUSH, POP, EXCHANGE) is pending, or needs to be recover. This can be concluded easily by the type of record stored in *Announce[pid]* (ExInfo or OpInfo), thus there is no need to explicitly know where exactly in the code the crash took place. Also, the RECOVER routine returns FAIL in case the last pending operation did not took affect (no linearization point), nor it will take in any future run. In such case, the user has the option to either re-invoke the operation, or to skip it, depends on the needs and circumstances of the specific use of the data structure.

The given implementation ignores the log of failures and successes of the exchange routine when recovering. That is, in case of a crash during an EXCHANGE, a process is able to recover the EXCHANGE routine, however, the log of successes and failures is not update, since it might be the process already updated it. In addition, in case of a FAIL response, we do not know whether the time limit (timeout) was reached, or that the process simply crashed earlier in the routine without completing it. The given implementation can be expanded to also consider the log. Nonetheless, for ease of presentation we do not handle the log in case of a crash. Assuming crash events are rare, the log still gives a roughly good approximation to the number of failures and successes, thus our approach might be useful in practice.

1.1 A Lock-Free Exchanger

An exchanger object supports the EXCHANGE procedure, which allows exactly two processes to exchange values. If process A calls the EXCHANGE with argument a , and B calls the EXCHANGE of the same object with argument b , then A's call will return value b and vice versa.

On the original algorithm [cite the book?!], processes race to win the exchanger using a **CAS** primitive. A process accessing the exchanger first reads its content, and act according to the state of it. The first process observe an EMPTY state, and tries to atomically writes its value and change the state to WAITING. In such case, it spins and wait for the second process to arrive. The second, observing the state is now WAITING, tries to write its value and change the state to BUSY. This way, it informs the first one a successful collision took place. Once the first process notice the collision, it reads the other process value and release the exchanger by setting it back to EMPTY. In order to avoid an unbounded waiting, if a second process does not show up, the call eventually timeout, and the process release the exchanger and return.

Assume a process p successfully capture the exchanger by setting its status to WAITING, followed by a crash. Now, some other process q complete the exchange by setting the exchanger to BUSY. Upon recovery, p can conclude some exchange was completed, but it can not tell whether its value is part of the exchange, and thus it can not complete the operation. Moreover, p and q must agree, otherwise q will return p 's value, and thus the operation of p must be linearized together with q operation.

In order to avoid the above problem, we take an approach resembling the BST implementation. Instead of writing a value to the exchanger, processes will use an info record, containing the relevant

information for the exchange. This way, processes use the exchanger in order to exchange info records (more precisely, pointers to such records), and not values. To overcome the problematic scenario described earlier, if a process q observe the exchanger state is WAITING with some record $yourop$, it first update its own record $myop$ it is about to try and collide with $yourop$, and only then performs the **CAS**. This way, if the collision is successful, the record $myop$ which now stored in the exchanger implies which two records collide. Also, the fact that different processes uses different records guarantee that at most one record can collide with $yourop$.

Using records instead of values, when using wisely, allows us to farther improve the algorithm. First, there is no need to store the exchanger's state in it (by using 2 bits of it to mark the state), but we can rather have this info in the record. Second, if there is a BUSY record in the exchanger, it contains the info of the two colliding records. Therefore, a third process, trying to also use the exchanger, can help the processes to complete the collision, and then can try and set the exchanger back to EMPTY, so it can use it again. In the original implementation, a process observing a BUSY exchanger, have to wait for the first process to read the value and release the exchanger. Therefore, if the first process crash after the collision, the exchanger will be hold by it forever. The helping mechanism avoids this scenario, making the exchange routine non-blocking.

Notice that no exchange record with EMPTY state is ever created, except for the *default* record. Therefore, reading EMPTY state is equivalent to the exchanger storing a pointer to *default*. A process p creates a new record $myop$ when accessing the exchanger, with a unique address. As long as p fails to perform a successful **CAS**, and thus fails to store $myop$ in *slot*, it is allowed to try again. However, once a process performs a successful **CAS** and stores $myop$ in *slot*, the only other **CAS** it is allowed to do are in order to try and store *default* in *slot*. Thus, $myop$ can be written exactly once to *slot*. It follows that a collision can occur between two processes exactly - once a WAITING record stored in *slot*, only a single **CAS** can replace it with a BUSY record. As the two records can not be written again to *slot*, no other process can collide with any of the records.

The EXCHANGE-RECOVER routine relies on the following argument. If a process p successfully wrote op_p to *slot* using the **CAS** in line 13, the only way to overwrite it by a different process q , is by a **CAS** in line 36 with a record op_q such that its state is BUSY, and $op_q.partner = op_p$. In addition, the only way to overwrite op_q is by a **CAS** replacing it with *default*, and this is done only after SWITCHPAIR(op_p, op_q) is completed, and thus both *result* fields are updated.

The correctness of the EXCHANGE-RECOVER routine is based on the above argument. There are few scenarios to consider. If p crash after a successful **CAS** in line 13, then op_p state is WAITING. Therefore, when reading *slot* in the EXCHANGE-RECOVER one of the following must hold. If *slot* contains op_p , then no process collide with p , and p continue to run as if the time limit has been reached. Otherwise, there was a collision. From the above argument, it must be that either op_q that collide with op_p is stored in *slot*, in this case $op_q.partner = op_p$, and p will try to complete the collision and release *slot*, or that op_q has been overwritten, and in this case the *result* field of op_p is updated. In both cases, p returns $op_p.result$. If p crash after a successful **CAS** in line 36, then op_p state is BUSY. It follows from the argument that the only way to overwrite op_p is only after completing the collision by SWITCHPAIR. Thus, either upon recovery p reads op_p from *slot*, and in this case it tries to complete the the operation, or that $op_p.result$ was already updated. In both cases, p returns it. If non of the above holds, then op_p was not involved in any collision, because either no successful **CAS** was done by p , or p reached the time limit while no process show up, and was able to set *slot* back to *default*. In any case, after the crash of p , op_p will never be written again to *slot*, nor any other op_q such that $op_q.partner = op_p$, as any such op_q tries to perform

CAS (op_p, op_q) that will fail. Also, as no process can collide with op_p , no SWITCHPAIR with op_p as parameter is ever invoked, and in particular $op_p.result = \perp$ for the rest of the execution. This in turn implies that upon recovery p will return FAIL, as required.

```

Type Node {
    T value
    int popby
    Node *next
}

Type PushInfo {                                ▷ subtype of Info
    Node *pushnd
}

Type PopInfo {                                  ▷ subtype of Info
    Node *popnd
}

Type ExInfo {                                   ▷ subtype of Info
    {EMPTY, WAITING, BUSY} state
    T value, result
    ExInfo *partner, *slot
}

```

Figure 1: Type definition

ExInfo *default* - global static ExInfo object with state = EMPTY

Algorithm 1: T EXCHANGE (ExInfo *slot, T *myitem*, long *timeout*)

```

1 long timeBound := getNanos() + timeout
2 ExInfo myop := new ExInfo(WAITING, myitem,  $\perp$ ,  $\perp$ , slot)
3 Announce[pid] := myop
4 while true do
5   if getNanos() > timeBound then
6     | myop.result := TIMEOUT // time limit reached
7     | return TIMEOUT
8   yourop := slot
9   switch yourop.state do
10    case EMPTY do
11      | myop.state := WAITING // attempt to replace default
12      | myop.partner :=  $\perp$ 
13      | if slot.CAS(yourop, myop) then // try to collide
14        | while getNanos() < timeBound do
15          | | yourop := slot
16          | | if yourop  $\neq$  myop then // a collision was done
17            | | | if yourop.parnter = myop then // yourop collide with myop
18              | | | SWITCHPAIR(myop, yourop)
19              | | | slot.CAS(yourop, default) // release slot
20            | | return myop.result
21          | end
22          | // time limit reached and no process collide with me
23          | if slot.CAS(myop, default) then // try to release slot
24            | | myop.result := TIMEOUT
25            | | return TIMEOUT
26          | else // some process show up
27            | | yourop := slot
28            | | if yourop.partner = myop then
29              | | | SWITCHPAIR(myop, yourop) // complete the collision
30              | | | slot.CAS(yourop, default) // release slot
31            | | return myop.result
32          | end
33        | break
34    case WAITING do // some process is waiting in slot
35      | myop.partner := yourop // attempt to replace yourop
36      | myop.state := BUSY
37      | if slot.CAS(yourop, myop) then // try to collide
38        | | SWITCHPAIR(myop, yourop) // complete the collision
39        | | slot.CAS(myop, default) // release slot
40        | | return myop.result
41      | break
42    case BUSY do // a collision in progress
43      | SWITCHPAIR(yourop, yourop.parnter) // help to complete the collision
44      | slot.CAS(yourop, default) // release slot
45      | break
46  end
47 end

```

Algorithm 2: void SWITCHPAIR(ExInfo <i>first</i> , ExInfo <i>second</i>)	
/* exchange the value of the two operations	*/
47 <i>first.result</i> := <i>second.value</i>	
48 <i>second.result</i> := <i>first.value</i>	
<hr/>	
Algorithm 3: T VISIT (T <i>value</i> , int <i>range</i> , long <i>duration</i>)	
/* invoke EXCHANGE on a random entry in the collision array	*/
49 int <i>cell</i> := randomNumber(<i>range</i>)	
50 return EXCHANGE(<i>exchanger</i> [<i>cell</i>], <i>value</i> , <i>duration</i>)	
<hr/>	
Algorithm 4: T EXCHANGE-RECOVER ()	
51 ExInfo * <i>myop</i> := Announce[<i>pid</i>]	// read your last operation record
52 ExInfo * <i>slot</i> := <i>myop.slot</i>	// and the slot on which it acts
53 if <i>myop.state</i> = WAITING then	
/* crash while trying to exchange <i>default</i> , or waiting for a collision	*/
54 <i>yourop</i> := <i>slot</i>	
55 if <i>yourop</i> = <i>myop</i> then	// still waiting for a collision
56 if <i>slot.CAS</i> (<i>myop</i> , <i>default</i>) then	// try to release slot
57 return FAIL	
58 else	// some process show up
59 <i>yourop</i> := <i>slot</i>	
60 if <i>yourop.partner</i> = <i>myop</i> then	
61 SWITCHPAIR(<i>myop</i> , <i>yourop</i>)	// complete the collision
62 <i>slot.CAS</i> (<i>yourop</i> , <i>default</i>)	// release slot
63 return <i>myop.result</i>	
64 else if <i>yourop.partner</i> = <i>myop</i> then	// <i>yourop</i> collide with <i>myop</i>
65 SWITCHPAIR(<i>myop</i> , <i>yourop</i>)	// complete the collision
66 <i>slot.CAS</i> (<i>yourop</i> , <i>default</i>)	// release slot
67 return <i>myop.result</i>	
68 if <i>myop.state</i> = BUSY then	
/* crash while trying to collide with <i>myop.partner</i>	*/
69 <i>yourop</i> := <i>slot</i>	
70 if <i>yourop</i> = <i>myop</i> then	// collide was successful and in progress
71 SWITCHPAIR(<i>myop</i> , <i>myop.partner</i>)	// complete the collision
72 <i>slot.CAS</i> (<i>myop</i> , <i>default</i>)	// release slot
73 return <i>myop.result</i>	
74 if <i>myop.result</i> ≠ ⊥ then	
75 return <i>myop.result</i>	// collide was successfully completed
76 else	
77 return FAIL	

Figure 2: Elimination Array routines

Algorithm 5: boolean TRY PUSH (Node *nd)

```
/* attempt to perform PUSH to the central stack */
78 Node *oldtop := Top
79 nd.next := oldtop
80 if Top.CAS(oldtop, nd) then           // try to declare nd as the new Head
81   | nd.popby.CAS( $\perp$ , NULL)           // announce nd is in the stack
82   | return true
83 return false
```

Algorithm 6: boolean PUSH (T myitem)

```
84 Node *nd = new Node (myitem)
85 nd.popby :=  $\perp$ 
86 PushInfo data := new PushInfo (nd)
87 while true do
88   | Announce[pid] := data           // declare - trying to push node nd
89   | if TRY PUSH(nd) then           // if central stack PUSH is successful
90     | return true
91   | range := CalculateRange()       // get parameters for collision array
92   | duration := CalculateDuration()
93   | othervalue := VISIT(myitem, range, duration)           // try to collide
94   | if othervalue = NULL then       // successfully collide with POP operation
95     | RecordSuccess ()
96     | return true
97   | else if othervalue = TIMEOUT then // failed to collide
98     | RecordFailure ()
99 end
```

Algorithm 7: boolean PUSH-ROCEOVER ()

```
100 Node *nd := Announce[pid].pushnd
101 if nd.popby  $\neq \perp$  then           // nd was announced to be in the stack
102   | return true
103 if SEARCH(nd) || nd.popby  $\neq \perp$  then // nd in the stack, or was announced as such
104   | nd.popby.CAS( $\perp$ , NULL)           // announce nd is in the stack
105   | return true
106 return FAIL
```

Algorithm 8: boolean SEARCH (Node *nd)

```
/* search for node nd in the stack */
107 Node *iter := Top
108 while iter  $\neq \perp$  do
109   | if iter = nd then
110     | return true
111   | iter := iter.next
112 end
113 return false
```

Figure 3: PUSH routine

Algorithm 9: T TRYPOP()

```
114 Node *oldtop := Top
115 Node *newtop
116 Announce[pid].popnd := oldnop           // declare - trying to pop node oldtop
117 if oldtop =  $\perp$  then                     // stack is empty
118 |   return EMPTY
119 newtop := oldtop.next
120 oldtop.popby.CAS( $\perp$ , NULL)               // announce oldtop is in the stack
121 if Top.CAS(oldtop, newtop) then // try to pop oldtop by changing Top to newtop
122 |   if newtop.popby.CAS(NULL, pid) then // try to announce yourself as winner
123 |   |   return oldtop.value
124 else
125 |   return  $\perp$ 
```

Algorithm 10: T POP ()

```
126 Node *result
127 PopInfo data := new PopInfo (Top)
128 while true do
129 |   Announce[pid] := data           // declare - trying to perform POP
130 |   result := TRYPOP()              // attempt to pop from central stack
131 |   if result  $\neq \perp$  then          // if central stack POP is successful
132 |   |   return result
133 |   range := CalculateRange()       // get parameters for collision array
134 |   duration := CalculateDuration()
135 |   othervalue := VISIT(NULL, range, duration) // try to collide
136 |   if othervalue = TIMEOUT then    // failed to collide
137 |   |   RecordFailure ()
138 |   else if othervalue  $\neq$  NULL then // successfully collide with PUSH operation
139 |   |   RecordSuccess ()
140 |   |   return othervalue
141 end
```

Algorithm 11: T POP-RECOVER()

```
142 Node *nd := Announce[pid].popnd // crash while trying to pop node nd
143 if nd =  $\perp$  then                  // pop from an empty stack
144 |   return EMPTY
145 if SEARCH(nd) then                // nd was not removed from the stack
146 |   return FAIL
147 nd.popby.CAS(NULL, pid)           // nd was removed. Try to complete the operation
148 if nd.popby = pid then             // you are the process to win the pop of nd
149 |   return nd.value
150 return FAIL
```

Figure 4: POP routine