

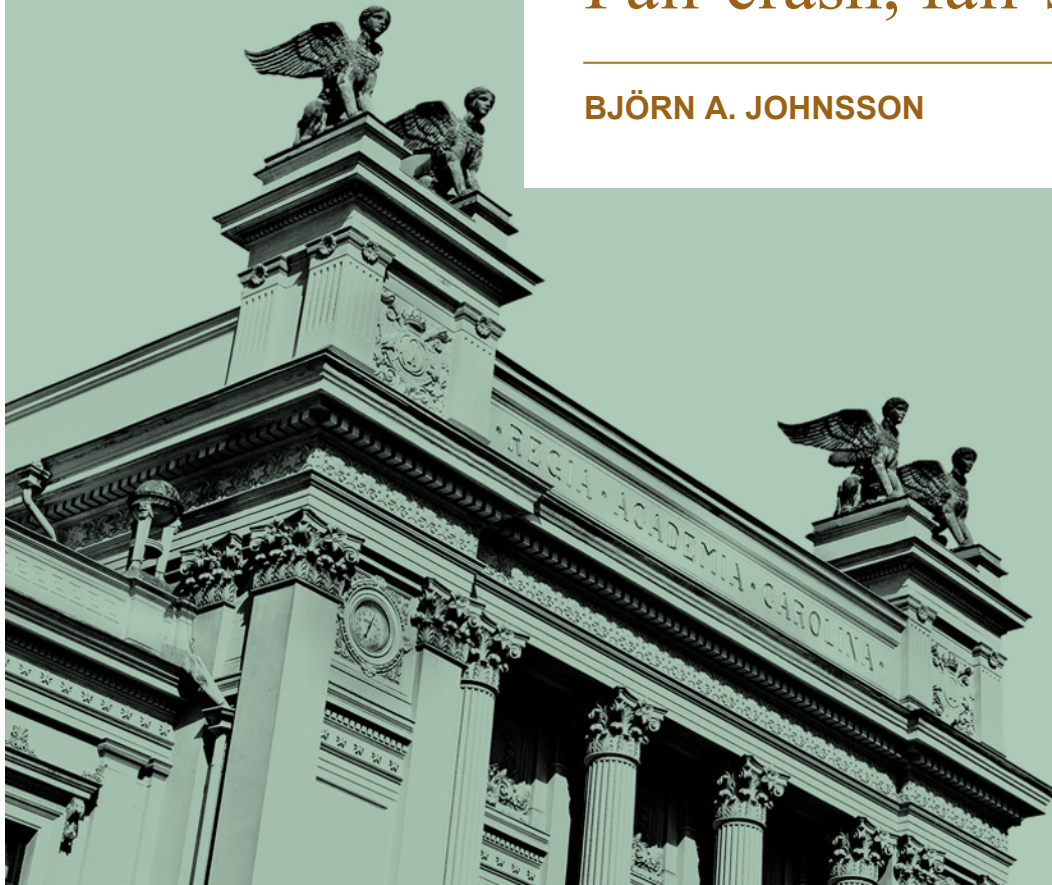


LUND
UNIVERSITY

Registers – Shared Memory

Fail-crash, fail-silent

BJÖRN A. JOHANSSON



Introduction

- Analogy from multi-CPU computers.
- Over network: *emulation of shared-memory*.
- Benefit: use shared-memory where there really is none.
- Considered easier than message exchanges.



Register Overview

- Process starts *read* operation with $\langle r, \textit{Read} \rangle$
- Process starts *write* operation with $\langle r, \textit{Write} \mid v \rangle$
- Process *completes* after reply event from register:
 - $\langle r, \textit{ReadReturn} \mid v \rangle$
 - $\langle r, \textit{WriteReturn} \rangle$
- Processes access registers in *sequential* manner
- Types: $(1, 1)$, $(1, N)$, (N, N)



Semantics

- *Liveness*: Every operation eventually completes.
- *Safety*: Every read operation returns the value written by the **last** write operation.

A process p invokes a write operation on a register with a value v and completes this write. Later on, some other process q invokes a write operation on the register with a new value w , and then q crashes before the operation completes. Hence, q does not get any indication that the operation has indeed taken place before it crashes, and the operation has failed. Now, if a process r subsequently invokes a read operation on the register, what is the value that r is supposed to return? Should it be v or w ?



Concurrency

- *Serial* (or *sequential*) exec: one operation after another
- *Concurrent* exec: what happens to def. of “last”?
- Three abstractions: *safe*, *regular*, and *atomic*.



Algorithm Overview

- $(1, N)$ Regular Register
 - Read-One Write-All
 - Majority Voting Regular Register
- $(1, N)$ Atomic Register
 - ~~$(1, N)$ Regular $\rightarrow (1, 1)$ Atomic $\rightarrow (1, N)$ Atomic Register~~
 - Read-Impose Write-All
 - Read-Impose Write-Majority
- (N, N) Atomic Register
 - ~~$(1, N)$ Atomic $\rightarrow (N, N)$ Atomic Register~~
 - Read-Impose Write-Consult-All
 - Read-Impose Write-Consult-Majority



Repetition

Distributed-System Models

- Fail-stop
 - crash-stop, perfekt links, perfect failure detector (P)
- Fail-silent
 - crash-stop, perfekt links, no failure detector



(1, N) Regular Register

Module 4.1: Interface and properties of a (1, N) regular register

Module:

Name: (1, N)-RegularRegister, **instance** *onrr*.

Events:

Request: $\langle onrr, Read \rangle$: Invokes a read operation on the register.

Request: $\langle onrr, Write \mid v \rangle$: Invokes a write operation with value v on the register.

Indication: $\langle onrr, ReadReturn \mid v \rangle$: Completes a read operation on the register with return value v .

Indication: $\langle onrr, WriteReturn \rangle$: Completes a write operation on the register.

Properties:

ONRR1: Termination: If a correct process invokes an operation, then the operation eventually completes.

ONRR2: Validity: A read that is not concurrent with a write returns the last value written; a read that is concurrent with a write returns the last value written or the value concurrently written.

initially \perp



(1, N) Regular Register

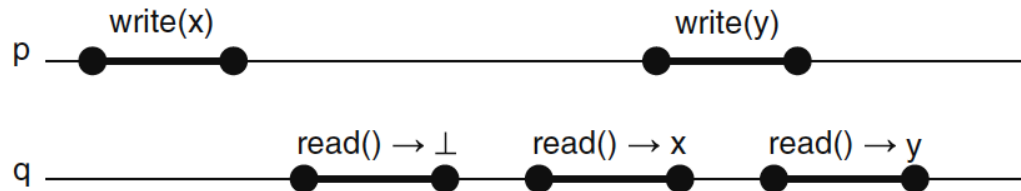
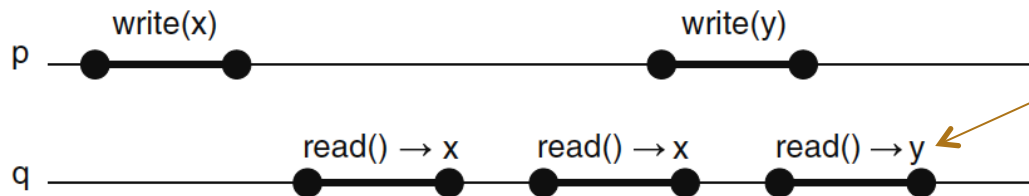


Figure 4.1: A register execution that is not regular because of the first read by process q



“x” also regular

Figure 4.2: A regular register execution



(1, N) Regular Register

Read-One Write-All

Algorithm 4.1: Read-One Write-All

Implements:

(1, N)-RegularRegister, **instance** *onrr*.

Uses:

BestEffortBroadcast, **instance** *beb*;
PerfectPointToPointLinks, **instance** *pl*;
PerfectFailureDetector, **instance** \mathcal{P} .

upon event $\langle \textit{onrr}, \textit{Init} \rangle$ **do**

val := \perp ;
correct := Π ;
writeset := \emptyset ;

upon event $\langle \mathcal{P}, \textit{Crash} \mid p \rangle$ **do**

correct := *correct* $\setminus \{p\}$;



(1, N) Regular Register

Read-One Write-All

```
upon event  $\langle onrr, Read \rangle$  do  
  trigger  $\langle onrr, ReadReturn \mid val \rangle$  ;  
  
upon event  $\langle onrr, Write \mid v \rangle$  do  
  trigger  $\langle \underline{beb}, Broadcast \mid [WRITE, v] \rangle$  ;  
  
upon event  $\langle beb, Deliver \mid q, [WRITE, v] \rangle$  do  
   $val := v$  ;  
  trigger  $\langle \underline{pl}, Send \mid q, ACK \rangle$  ;  
  
upon event  $\langle pl, Deliver \mid p, ACK \rangle$  do  
   $writeset := writeset \cup \{p\}$  ;  
  
upon correct  $\subseteq writeset$  do  
   $writeset := \emptyset$  ;  
  trigger  $\langle onrr, WriteReturn \rangle$  ;
```



(1, N) Regular Register

Read-One Write-All

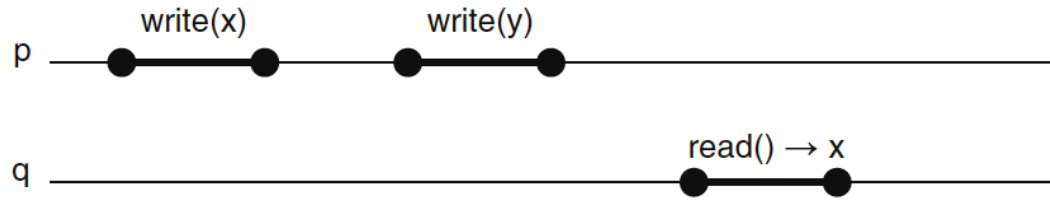


Figure 4.3: A non-regular register execution



(1, N) Regular Register

Majority Voting Regular Register

Algorithm 4.2: Majority Voting Regular Register

Implements:

(1, N)-RegularRegister, **instance** *onrr*.

Uses:

BestEffortBroadcast, **instance** *beb*;

PerfectPointToPointLinks, **instance** *pl*.

upon event $\langle onrr, Init \rangle$ do

$(ts, val) := (0, \perp)$;

$wts := 0$;

$acks := 0$;

$rid := 0$;

$readlist := [\perp]^N$;

upon event $\langle onrr, Write \mid v \rangle$ do

$wts := wts + 1$;

$acks := 0$;

trigger $\langle beb, Broadcast \mid [WRITE, wts, v] \rangle$;

upon event $\langle beb, Deliver \mid p, [WRITE, ts', v'] \rangle$ do

if $ts' > ts$ **then**

$(ts, val) := (ts', v')$;

trigger $\langle pl, Send \mid p, [ACK, ts'] \rangle$;



(1, N) Regular Register

Majority Voting Regular Register

```

upon event  $\langle pl, Deliver \mid q, [ACK, ts'] \rangle$  such that  $ts' = wts$  do
   $acks := acks + 1;$ 
  if  $acks > N/2$  then
     $acks := 0;$ 
    trigger  $\langle onrr, WriteReturn \rangle;$ 

```

```

upon event  $\langle onrr, Read \rangle$  do
   $rid := rid + 1;$ 
   $readlist := [\perp]^N;$ 
  trigger  $\langle beb, Broadcast \mid [READ, rid] \rangle;$ 

```

```

upon event  $\langle beb, Deliver \mid p, [READ, r] \rangle$  do
  trigger  $\langle pl, Send \mid p, [VALUE, r, ts, val] \rangle;$ 

```

```

upon event  $\langle pl, Deliver \mid q, [VALUE, r, ts', v'] \rangle$  such that  $r = rid$  do
   $readlist[q] := (ts', v');$ 
  if  $\#(readlist) > N/2$  then
     $v := \text{highestval}(readlist);$ 
     $readlist := [\perp]^N;$ 
    trigger  $\langle onrr, ReadReturn \mid v \rangle;$ 

```

← returns pair with greatest time stamp



(1, N) Atomic Register

Module 4.2: Interface and properties of a (1, N) atomic register

Module:

Name: (1, N)-AtomicRegister, **instance** *onar*.

Events:

Request: $\langle onar, Read \rangle$: Invokes a read operation on the register.

Request: $\langle onar, Write \mid v \rangle$: Invokes a write operation with value v on the register.

Indication: $\langle onar, ReadReturn \mid v \rangle$: Completes a read operation on the register with return value v .

Indication: $\langle onar, WriteReturn \rangle$: Completes a write operation on the register.

Properties:

ONAR1–ONAR2: Same as properties ONRR1–ONRR2 of a (1, N) regular register (Module 4.1).

! **ONAR3: Ordering:** If a read returns a value v and a subsequent read returns a value w , then the write of w does not precede the write of v .



(1, N) Atomic Register

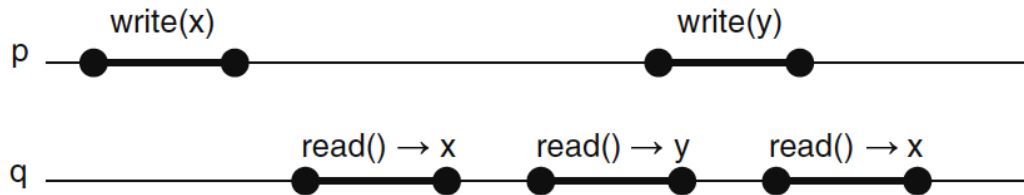


Figure 4.4: A register execution that is not atomic because of the third read by process q

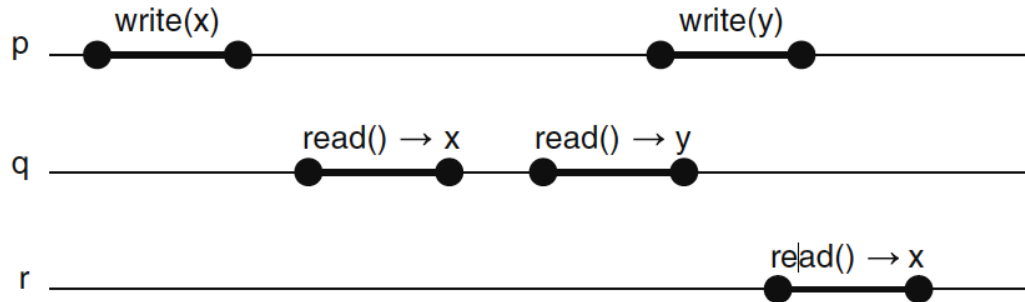


Figure 4.5: Violation of atomicity in the “Read-One Write-All” regular register algorithm



(1, N) Atomic Register

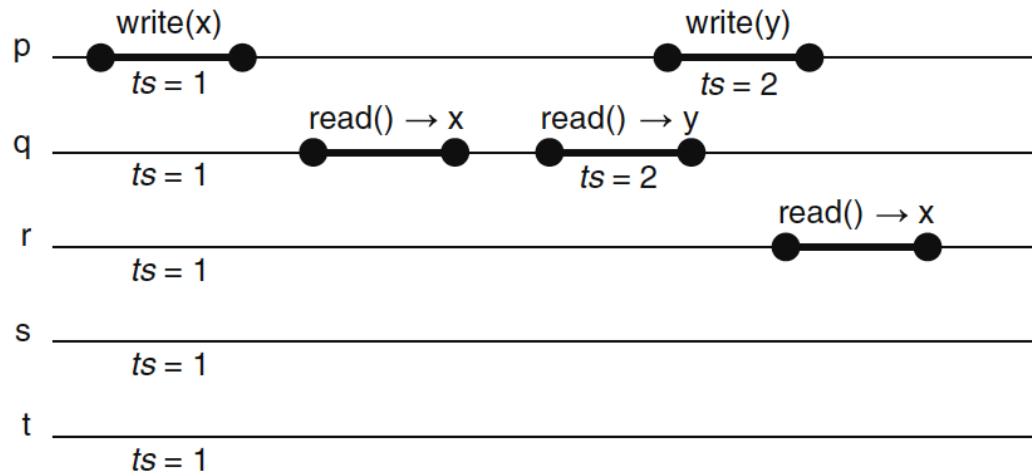


Figure 4.6: Violation of atomicity in the “Majority Voting” regular register algorithm



(1, N) Atomic Register

(1, N) Regular \rightarrow (1, 1) Atomic Register

Algorithm 4.3: From (1, N) Regular to (1, 1) Atomic Registers

Implements:

(1, 1)-AtomicRegister, **instance** *oarr*.

Uses:

(1, N)-RegularRegister, **instance** *onrr*.

upon event $\langle oarr, Init \rangle$ **do**

$(ts, val) := (0, \perp)$;

$wts := 0$;

upon event $\langle oarr, Write \mid v \rangle$ **do**

$wts := wts + 1$;

trigger $\langle onrr, Write \mid (wts, v) \rangle$;

upon event $\langle onrr, WriteReturn \rangle$ **do**

trigger $\langle oarr, WriteReturn \rangle$;

upon event $\langle oarr, Read \rangle$ **do**

trigger $\langle onrr, Read \rangle$;

upon event $\langle onrr, ReadReturn \mid (ts', v') \rangle$ **do**

if $ts' > ts$ **then**

$(ts, val) := (ts', v')$;

trigger $\langle oarr, ReadReturn \mid val \rangle$;



$(1, N)$ Atomic Register

$(1, 1)$ Atomic \rightarrow $(1, N)$ Atomic Register

Algorithm 4.4: From $(1, 1)$ Atomic to $(1, N)$ Atomic Registers

Implements:

$(1, N)$ -AtomicRegister, **instance** *onar*.

Uses:

$(1, 1)$ -AtomicRegister (multiple instances).

upon event $\langle onar, Init \rangle$ **do**

$ts := 0;$

$acks := 0;$

$writing := \text{FALSE};$

$readval := \perp;$

$readlist := [\perp]^N;$

forall $q \in \Pi, r \in \Pi$ **do**

Initialize a new instance $ooar.q.r$ of $(1, 1)$ -AtomicRegister
with writer r and reader q ;

upon event $\langle onar, Write \mid v \rangle$ **do**

$ts := ts + 1;$

$writing := \text{TRUE};$

forall $q \in \Pi$ **do**

trigger $\langle ooar.q.self, Write \mid (ts, v) \rangle;$



(1, N) Atomic Register

(1, 1) Atomic \rightarrow (1, N) Atomic Register

```

upon event  $\langle ooar.q.self, WriteReturn \rangle$  do
   $acks := acks + 1;$ 
  if  $acks = N$  then
     $acks := 0;$ 
    if  $writing = TRUE$  then
      trigger  $\langle onar, WriteReturn \rangle;$ 
       $writing := FALSE;$ 
    else
      trigger  $\langle onar, ReadReturn \mid readval \rangle;$ 

```

```

upon event  $\langle onar, Read \rangle$  do
  forall  $r \in \Pi$  do
    trigger  $\langle ooar.self.r, Read \rangle;$ 

```

```

upon event  $\langle ooar.self.r, ReadReturn \mid (ts', v') \rangle$  do
   $readlist[r] := (ts', v');$ 
  if  $\#(readlist) = N$  then
     $(maxts, readval) := highest(readlist);$ 
     $readlist := [\perp]^N;$ 
    forall  $q \in \Pi$  do
      trigger  $\langle ooar.q.self, Write \mid (maxts, readval) \rangle;$ 

```



(1, N) Atomic Register

Read-Impose Write-All

Algorithm 4.5: Read-Impose Write-All

Implements:

(1, N)-AtomicRegister, **instance** *onar*.

Uses:

BestEffortBroadcast, **instance** *beb*;

PerfectPointToPointLinks, **instance** *pl*;

PerfectFailureDetector, **instance** \mathcal{P} .

upon event $\langle onar, Init \rangle$ **do**

$(ts, val) := (0, \perp)$;

$correct := \Pi$;

$writeset := \emptyset$;

$readval := \perp$;

$reading := \text{FALSE}$;

upon event $\langle \mathcal{P}, Crash \mid p \rangle$ **do**

$correct := correct \setminus \{p\}$;

upon event $\langle onar, Read \rangle$ **do**

! $reading := \text{TRUE}$;

$readval := val$;

trigger $\langle beb, Broadcast \mid [\text{WRITE}, ts, val] \rangle$;



same for $\langle onar, Write \rangle$!



(1, N) Atomic Register

Read-Impose Write-All

```

upon event  $\langle onar, Write \mid v \rangle$  do
  trigger  $\langle beb, Broadcast \mid [WRITE, ts + 1, v] \rangle$ ;

```

```

upon event  $\langle beb, Deliver \mid p, [WRITE, ts', v'] \rangle$  do
  if  $ts' > ts$  then
     $(ts, val) := (ts', v')$ ;
  trigger  $\langle pl, Send \mid p, [ACK] \rangle$ ;

```

same for $\langle onar, Read \rangle$!

```

upon event  $\langle pl, Deliver \mid p, [ACK] \rangle$  then
   $writeset := writeset \cup \{p\}$ ;

```

```

upon correct  $\subseteq writeset$  do
   $writeset := \emptyset$ ;
  if  $reading = TRUE$  then
     $reading := FALSE$ ;
    trigger  $\langle onar, ReadReturn \mid readval \rangle$ ;
  else
    trigger  $\langle onar, WriteReturn \rangle$ ;

```



(1, N) Atomic Register

Read-Improve Write-Majority

Algorithm 4.6: Read-Improve Write-Majority (part 1, read)

Implements:

(1, N)-AtomicRegister, **instance** *onar*.

Uses:

BestEffortBroadcast, **instance** *beb*;

PerfectPointToPointLinks, **instance** *pl*.

upon event $\langle onar, Init \rangle$ **do**

$(ts, val) := (0, \perp)$;

$wts := 0$;

$acks := 0$;

$rid := 0$;

$readlist := [\perp]^N$;

$readval := \perp$;

$reading := \text{FALSE}$;



(1, N) Atomic Register

Read-Impose Write-Majority

upon event $\langle onar, Read \rangle$ **do**

$rid := rid + 1;$

$acks := 0;$

$readlist := [\perp]^N;$

$reading := TRUE;$

trigger $\langle beb, Broadcast \mid [READ, rid] \rangle;$

upon event $\langle beb, Deliver \mid p, [READ, r] \rangle$ **do**

trigger $\langle pl, Send \mid p, [VALUE, r, ts, val] \rangle;$

upon event $\langle pl, Deliver \mid q, [VALUE, r, ts', v'] \rangle$ **such that** $r = rid$ **do**

$readlist[q] := (ts', v');$

if $\#(readlist) > N/2$ **then**

$(maxts, readval) := highest(readlist);$

$readlist := [\perp]^N;$

trigger $\langle beb, Broadcast \mid [WRITE, rid, maxts, readval] \rangle;$



(1, N) Atomic Register

Read-Improve Write-Majority

```

upon event  $\langle onar, Write \mid v \rangle$  do
   $rid := rid + 1;$ 
   $wts := wts + 1;$ 
   $acks := 0;$ 
  trigger  $\langle beb, Broadcast \mid [WRITE, rid, wts, v] \rangle;$ 

upon event  $\langle beb, Deliver \mid p, [WRITE, r, ts', v'] \rangle$  do
  if  $ts' > ts$  then
     $(ts, val) := (ts', v');$ 
  trigger  $\langle pl, Send \mid p, [ACK, r] \rangle;$ 

upon event  $\langle pl, Deliver \mid q, [ACK, r] \rangle$  such that  $r = rid$  do
   $acks := acks + 1;$ 
  if  $acks > N/2$  then
     $acks := 0;$ 
    if  $reading = TRUE$  then
       $reading := FALSE;$ 
      trigger  $\langle onar, ReadReturn \mid readval \rangle;$ 
    else
      trigger  $\langle onar, WriteReturn \rangle;$ 

```



(N, N) Atomic Register

Module 4.3: Interface and properties of an (N, N) atomic register

Module:

Name: (N, N) -AtomicRegister, **instance** *nmar*.

Events:

Request: $\langle nmar, Read \rangle$: Invokes a read operation on the register.

Request: $\langle nmar, Write \mid v \rangle$: Invokes a write operation with value v on the register.

Indication: $\langle nmar, ReadReturn \mid v \rangle$: Completes a read operation on the register with return value v .

Indication: $\langle nmar, WriteReturn \rangle$: Completes a write operation on the register.

Properties:

NNAR1: *Termination:* Same as property ONAR1 of a $(1, N)$ atomic register (Module 4.2).

! **NNAR2:** *Atomicity:* Every read operation returns the value that was written most recently in a hypothetical execution, where every failed operation appears to be complete or does not appear to have been invoked at all, and every complete operation appears to have been executed at some instant between its invocation and its completion.



(N, N) Atomic Register

The hypothetical serial execution mentioned before is called a *linearization* of the actual execution. More precisely, a linearization of an execution is defined as a sequence of complete operations that appear atomically, one after the other, which contains at least all complete operations of the actual execution (and possibly some operations that were incomplete) and satisfies the following conditions:

1. every read returns the last value written; and
2. for any two operations o and o' , if o precedes o' in the actual execution, then o also appears before o' in the linearization.

We call an execution *linearizable* if there is a way to *linearize* it like this. With this notion, one can reformulate the *atomicity* property of an (N, N) atomic register in Module 4.3 as:

NNAR2': Atomicity: Every execution of the register is linearizable.



(N, N) Atomic Register

$(1, N)$ Atomic \rightarrow (N, N) Atomic Register

Algorithm 4.8: From $(1, N)$ Atomic to (N, N) Atomic Registers

Implements:

(N, N) -AtomicRegister, **instance** *n nar*.

Uses:

$(1, N)$ -AtomicRegister (multiple instances).

upon event $\langle n nar, Init \rangle$ **do**

val := \perp ;

writing := FALSE;

readlist := $[\perp]^N$;

forall $q \in \Pi$ **do**

Initialize a new instance *onar.q* of $(1, N)$ -AtomicRegister
with writer *q*;



(N, N) Atomic Register

$(1, N)$ Atomic \rightarrow (N, N) Atomic Register

```

upon event  $\langle nnar, Write \mid v \rangle$  do
   $val := v;$ 
   $writing := TRUE;$ 
  forall  $q \in \Pi$  do
    trigger  $\langle onar.q, Read \rangle;$ 

upon event  $\langle nnar, Read \rangle$  do
  forall  $q \in \Pi$  do
    trigger  $\langle onar.q, Read \rangle;$ 

upon event  $\langle onar.q, ReadReturn \mid (ts', v') \rangle$  do
   $readlist[q] := (ts', rank(q), v');$ 
  if  $\#(readlist) = N$  then
     $(ts, v) := highest(readlist);$ 
     $readlist := [\perp]^N;$ 
    if  $writing = TRUE$  then
       $writing := FALSE;$ 
      trigger  $\langle onar.self, Write \mid (ts + 1, val) \rangle;$ 
    else
      trigger  $\langle nnar, ReadReturn \mid v \rangle;$ 

upon event  $\langle onar.self, WriteReturn \rangle$  do
  trigger  $\langle nnar, WriteReturn \rangle;$ 

```



(N, N) Atomic Register

Read-Improve Write-Consult-All

Algorithm 4.9: Read-Improve Write-Consult-All

Implements: (N, N) -AtomicRegister, **instance** *nmar*.**Uses:**BestEffortBroadcast, **instance** *beb*;
PerfectPointToPointLinks, **instance** *pl*;
PerfectFailureDetector, **instance** \mathcal{P} .**upon event** $\langle nmar, Init \rangle$ **do****!** $(ts, \underline{wr}, val) := (0, 0, \perp)$;
 $correct := \Pi$;
 $writeset := \emptyset$;
 $readval := \perp$;
 $reading := \text{FALSE}$;**upon event** $\langle \mathcal{P}, Crash \mid p \rangle$ **do** $correct := correct \setminus \{p\}$;

(N, N) Atomic Register

Read-Impose Write-Consult-All

```

upon event  $\langle nnar, Read \rangle$  do
  reading := TRUE;
  readval := val;
  trigger  $\langle beb, Broadcast \mid [WRITE, ts, wr, val] \rangle$ ;

upon event  $\langle nnar, Write \mid v \rangle$  do
  trigger  $\langle beb, Broadcast \mid [WRITE, ts + 1, rank(self), v] \rangle$ ;

upon event  $\langle beb, Deliver \mid p, [WRITE, ts', wr', v'] \rangle$  do
  if  $(ts', wr')$  is larger than  $(ts, wr)$  then
     $(ts, wr, val) := (ts', wr', v')$ ;
  trigger  $\langle pl, Send \mid p, [ACK] \rangle$ ;

upon event  $\langle pl, Deliver \mid p, [ACK] \rangle$  then
  writeset := writeset  $\cup \{p\}$ ;

upon correct  $\subseteq$  writeset do
  writeset :=  $\emptyset$ ;
  if reading = TRUE then
    reading := FALSE;
    trigger  $\langle nnar, ReadReturn \mid readval \rangle$ ;
  else
    trigger  $\langle nnar, WriteReturn \rangle$ ;

```



(N, N) Atomic Register

Read-Improve Write-Consult-Majority

Algorithm 4.10: Read-Improve Write-Consult-Majority (part 1, read and consult)

Implements:

(N, N) -AtomicRegister, **instance** *nmar*.

Uses:

BestEffortBroadcast, **instance** *beb*;

PerfectPointToPointLinks, **instance** *pl*.

upon event $\langle nmar, Init \rangle$ **do**

! $(ts, \underline{wr}, val) := (0, 0, \perp)$;

$acks := 0$;

$\underline{writeval} := \perp$;

$rid := 0$;

$readlist := [\perp]^N$;

$readval := \perp$;

$reading := \text{FALSE}$;



(N, N) Atomic Register

Read-Impose Write-Consult-Majority

upon event $\langle nna, Read \rangle$ **do**

$rid := rid + 1;$

$acks := 0;$

$readlist := [\perp]^N;$

$reading := TRUE;$

trigger $\langle beb, Broadcast \mid [READ, rid] \rangle;$

upon event $\langle beb, Deliver \mid p, [READ, r] \rangle$ **do**

trigger $\langle pl, Send \mid p, [VALUE, r, ts, wr, val] \rangle;$

upon event $\langle pl, Deliver \mid q, [VALUE, r, ts', wr', v'] \rangle$ **such that** $r = rid$ **do**

$readlist[q] := (ts', wr', v');$

if $\#(readlist) > N/2$ **then**

$(maxts, rr, readval) := highest(readlist);$

$readlist := [\perp]^N;$

if $reading = TRUE$ **then**

trigger $\langle beb, Broadcast \mid [WRITE, rid, maxts, rr, readval] \rangle;$

else

trigger $\langle beb, Broadcast \mid [WRITE, rid, maxts + 1, rank(self), writeval] \rangle;$

factors in *rank*



(N, N) Atomic Register

Read-Impose Write-Consult-Majority

```

upon event  $\langle nnar, Write \mid v \rangle$  do
   $rid := rid + 1;$ 
   $writeval := v;$ 
   $acks := 0;$ 
   $readlist := [\perp]^N;$ 
  trigger  $\langle beb, Broadcast \mid [READ, rid] \rangle;$ 

upon event  $\langle beb, Deliver \mid p, [WRITE, r, ts', wr', v'] \rangle$  do
  if  $(ts', wr')$  is larger than  $(ts, wr)$  then
     $(ts, wr, val) := (ts', wr', v');$ 
  trigger  $\langle pl, Send \mid p, [ACK, r] \rangle;$ 

upon event  $\langle pl, Deliver \mid q, [ACK, r] \rangle$  such that  $r = rid$  do
   $acks := acks + 1;$ 
  if  $acks > N/2$  then
     $acks := 0;$ 
    if  $reading = TRUE$  then
       $reading := FALSE;$ 
      trigger  $\langle nnar, ReadReturn \mid readval \rangle;$ 
    else
      trigger  $\langle nnar, WriteReturn \rangle;$ 

```





LUND
UNIVERSITY