

Atomic register algorithms

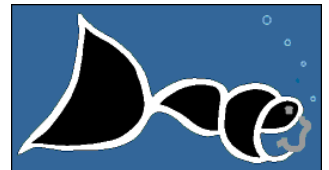
R. Guerraoui

*Distributed Programming Laboratory
lpdwww.epfl.ch*



© R. Guerraoui

1



Overview of this lecture

- ***(1) From regular to atomic***
- ***(2) A 1-1 atomic fail-stop algorithm***
- ***(3) A 1-N atomic fail-stop algorithm***
- ***(4) A N-N atomic fail-stop algorithm***
- ***(5) From fail-stop to fail-silent***

Fail-stop algorithms

- We first assume a fail-stop model; more precisely:
 - any number of processes can fail by crashing (no recovery)
 - channels are reliable
 - failure detection is perfect

The simple algorithm

- Consider our fail-stop ***regular*** register algorithm
 - every process has a local copy of the register value
 - every process reads ***locally***
 - the writer writes ***globally***, i.e., at all (non-crashed) processes

The simple algorithm

Write(v) at p_i

- send $[W, v]$ to all
- for every p_j , wait until either:
 - received $[\text{ack}]$ or
 - suspected $[p_j]$
- Return ok

At p_i :

when receive $[W, v]$
from p_j

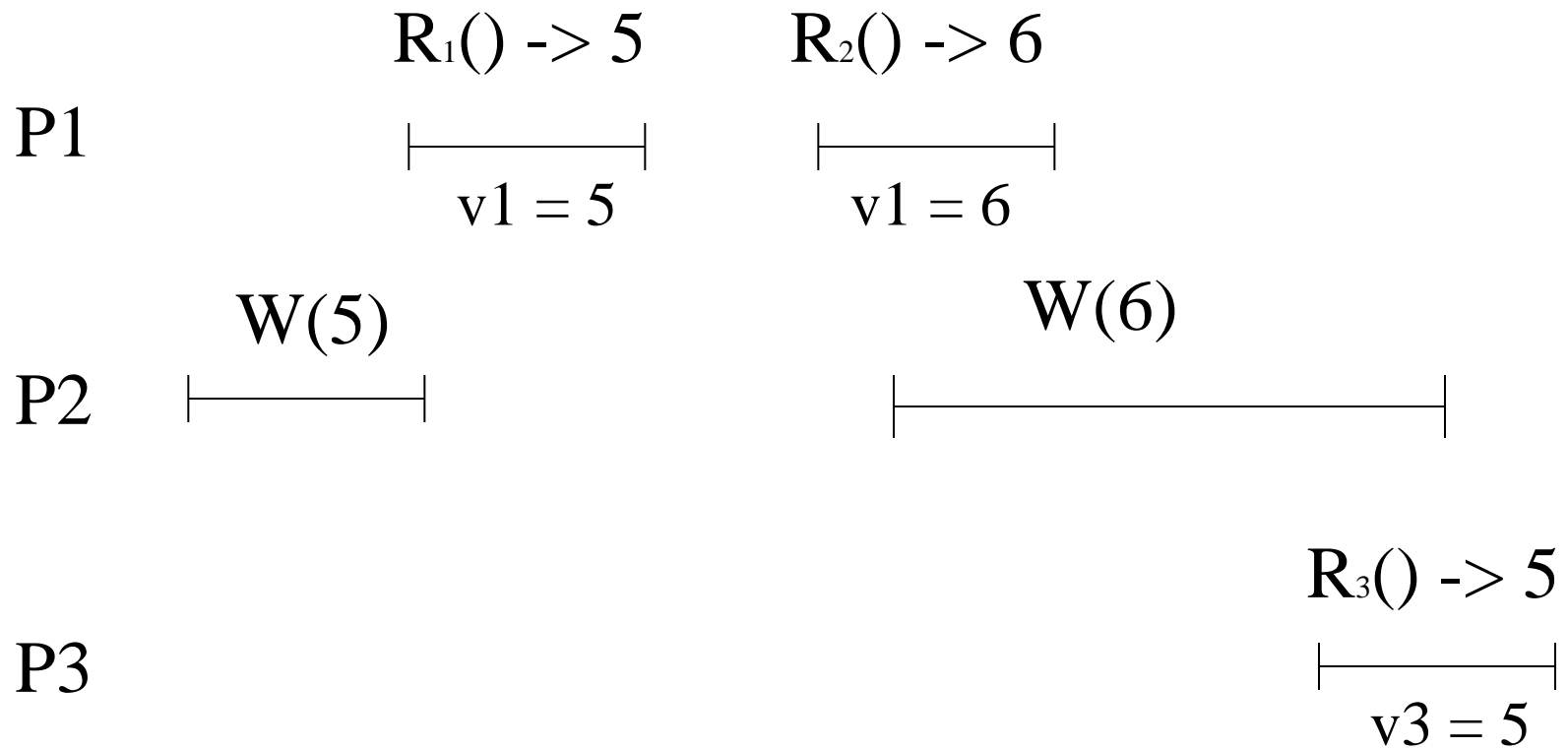
$v_i := v$

send $[\text{ack}]$ to p_j

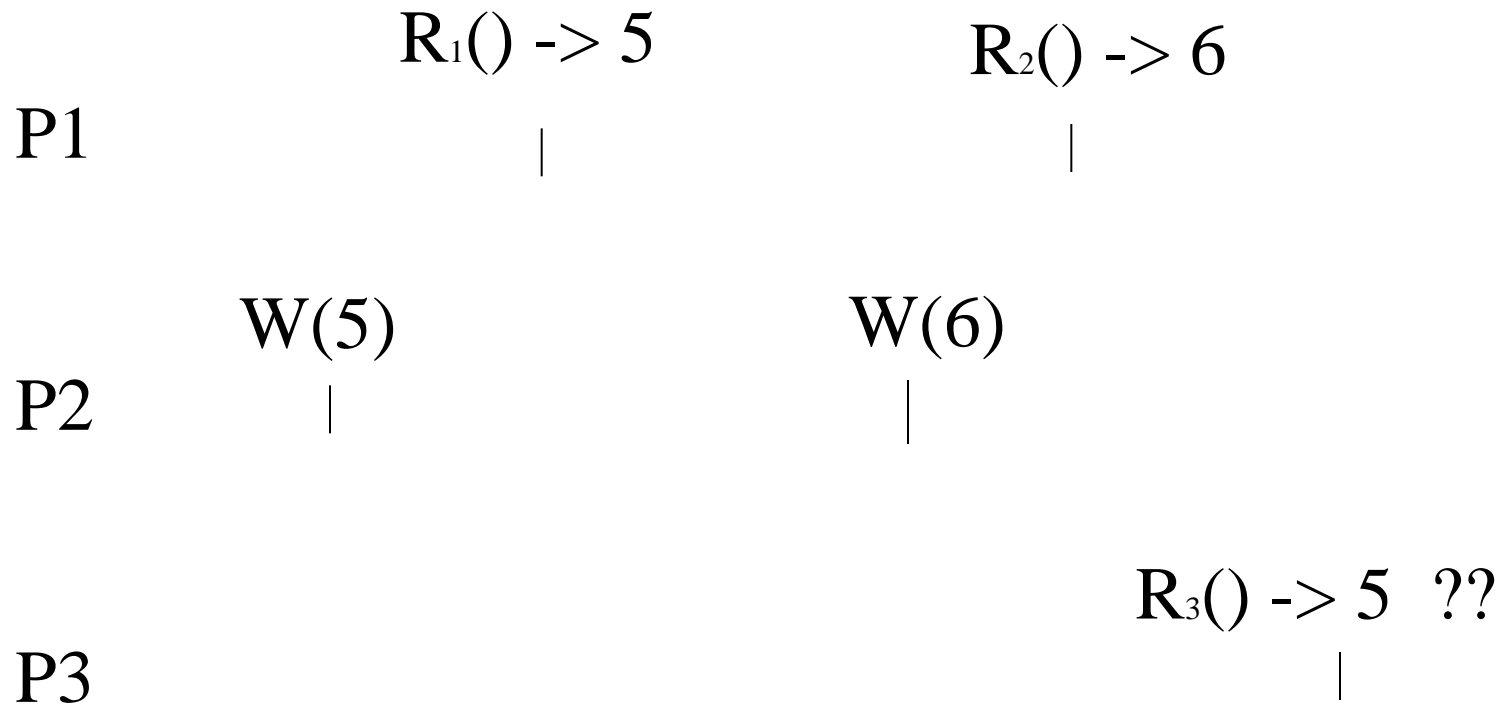
Read() at p_i

- Return v_i

Atomicity?



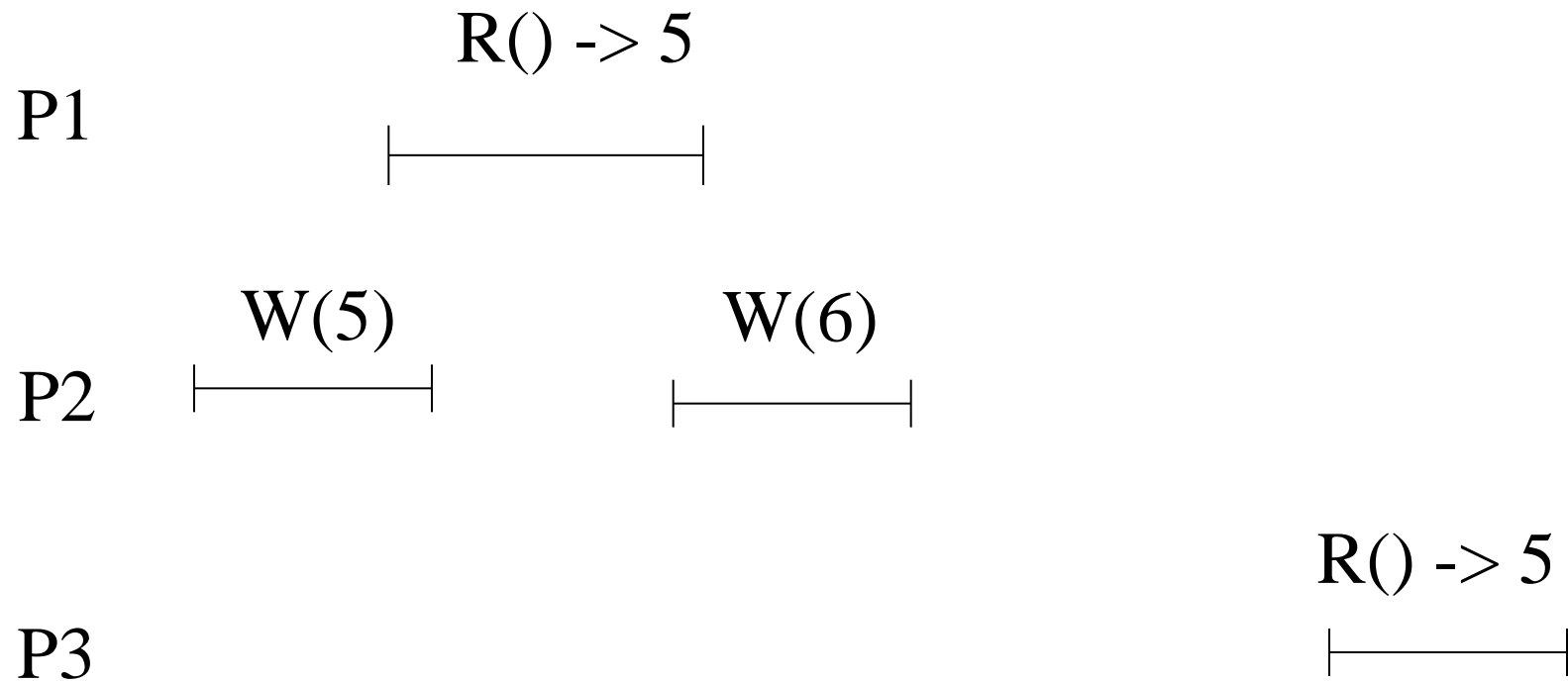
Linearization?



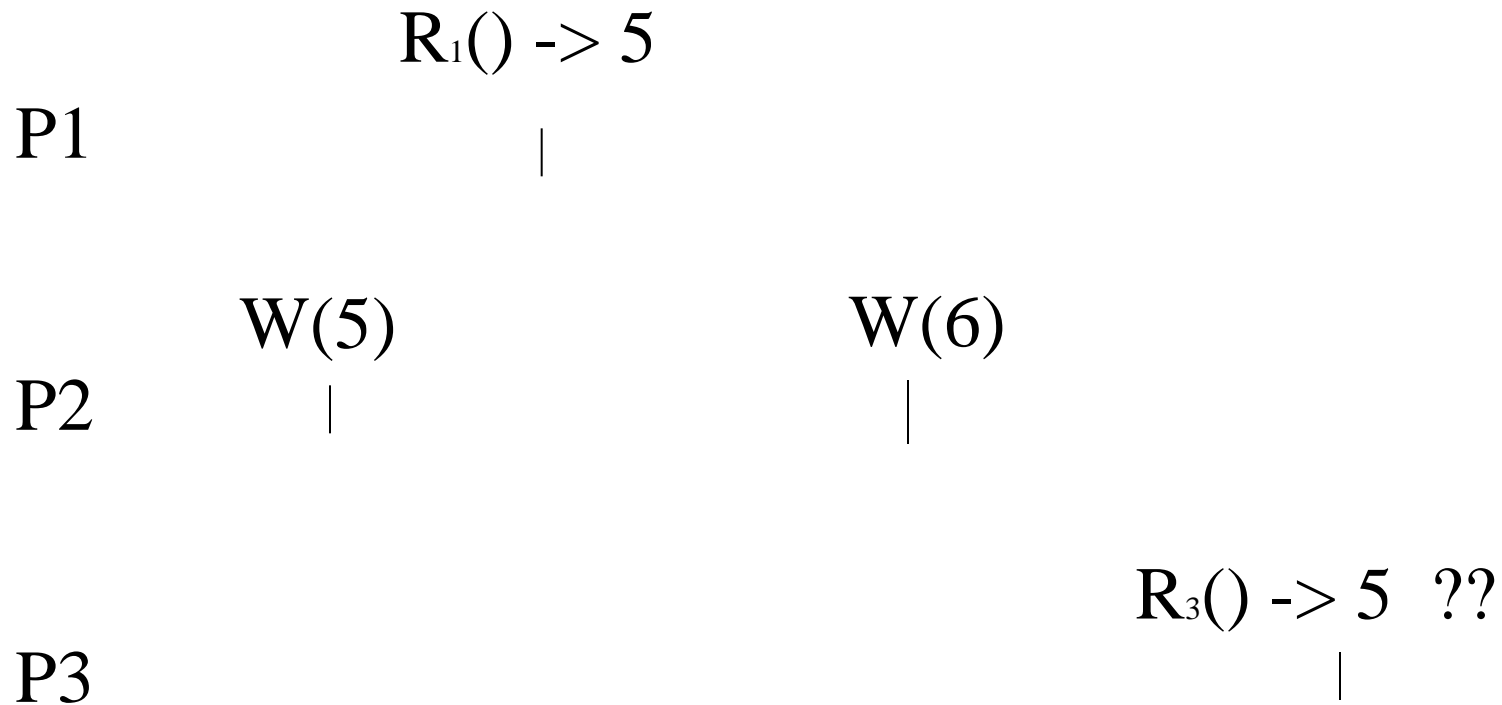
Fixing the pb: read-globally

- Read() at p_i
 - send $[W, v_i]$ to all
 - for every p_j , wait until either:
 - receive $[ack]$ or
 - suspect $[p_j]$
 - Return v_i

Still a problem



Linearization?



Overview of this lecture

- ***(1) From regular to atomic***
- ***(2) A 1-1 atomic fail-stop algorithm***
- ***(3) A 1-N atomic fail-stop algorithm***
- ***(4) A N-N atomic fail-stop algorithm***
- ***(5) From fail-stop to fail-silent***

A fail-stop 1-1 atomic algorithm

Write(v) at p_1

- send $[W, v]$ to p_2
- Wait until either:
 - receive $[ack]$ from p_2 or
 - suspect $[p_2]$
- Return ok

At p_2 :

when receive $[W, v]$ from p_1
 $v_2 := v$
send $[ack]$ to p_2

Read() at p_2

- Return v_2

A fail-stop 1-N algorithm

- every process maintains a local value of the register as well as a sequence number
- the writer, p_1 , maintains, in addition a timestamp ts_1
- any process can read in the register

A fail-stop 1-N algorithm

Write(v) at p_1

- $ts1++$
- send $[W,ts1,v]$ to all
- for every p_i , wait until either:
 - receive $[ack]$ or
 - suspect $[p_i]$
- Return ok

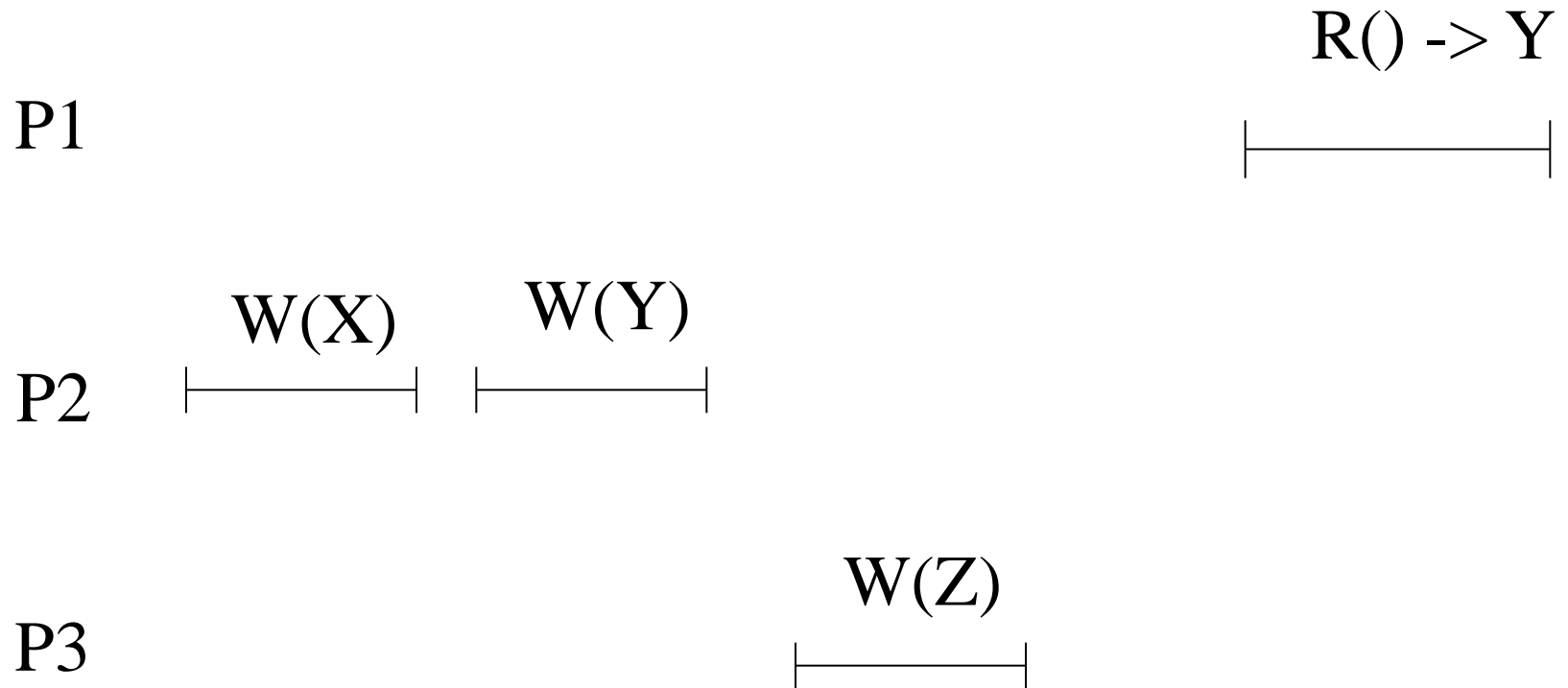
Read() at p_i

- send $[W,sni,vi]$ to all
- for every p_j , wait until either:
 - receive $[ack]$ or
 - suspect $[p_j]$
- Return vi

A 1-N algorithm (cont'd)

- At p_i
 - When p_i receive $[W, ts, v]$ from p_j
 - if $ts > sn_i$ then
 - $v_i := v$
 - $sn_i := ts$
 - send $[ack]$ to p_j

Why not N-N?



The Write() algorithm

- Write(v) at p_i
 - ✓ send $[W]$ to all
 - ✓ for every p_j wait until
 - **receive $[W,sn_j]$ or**
 - **suspect p_j**
 - ✓ $(sn,id) := (\text{highest } sn_j + 1,i)$
 - ✓ send $[W,(sn,id),v]$ to all
 - ✓ for every p_j wait until
 - **receive $[W,(sn,id),ack]$ or**
 - **suspect p_j**
 - ✓ Return ok
- At p_i
 - T1:
 - ✓ when receive $[W]$ from p_j
 - send $[W,sn]$ to p_j
 - T2:
 - ✓ when receive $[W,(sn_j,id_j),v]$ from p_j
 - ✓ If $(sn_j,id_j) > (sn,id)$ then
 - $v_i := v$
 - $(sn,id) := (sn_j,id_j)$
 - ✓ send $[W,(sn_j,id_j),ack]$ to p_j

The Read() algorithm

- Read() at p_i
 - ✓ send [R] to all
 - ✓ for every p_j wait until
 - **receive [R,(sn_j,id_j),v_j] or**
 - **suspect p_j**
 - ✓ $v = v_j$ with the highest (sn_j,id_j)
 - ✓ (sn,id) = highest (sn_j,id_j)
 - ✓ send [W,(sn,id),v] to all
 - ✓ for every p_j wait until
 - **receive [W,(sn,id),ack] or**
 - **suspect p_j**
 - ✓ Return v
- At p_i
 - T1:
 - ✓ when receive [R] from p_j
 - send [R,(sn,id),v_i] to p_j
 - T2:
 - ✓ when receive [W,(sn_j,id_j),v] from p_j
 - ✓ If (sn_j,id_j) > (sn,id) then
 - $v_i := v$
 - (sn,id) := (sn_j,id_j)
 - ✓ send [W,(sn_j,id_j),ack] to p_j

Overview of this lecture

- ***(1) From regular to atomic***
- ***(2) A 1-1 atomic fail-stop algorithm***
- ***(3) A 1-N atomic fail-stop algorithm***
- ***(4) A N-N atomic fail-stop algorithm***
- ***(5) From fail-stop to fail-silent***

From fail-stop to fail-silent

- We assume a majority of correct processes
- In the 1-N algorithm, the writer writes in a majority using a timestamp determined locally and the reader selects a value from a majority and then imposes this value on a majority
- In the N-N algorithm, the writers determines first the timestamp using a majority