

Implementing the Consensus Object with Timing Assumptions

R. Guerraoui
Distributed Programming Laboratory



© R. Guerraoui

1



A modular approach

We implement *Wait-free Consensus (Consensus)*
through:

Lock-free Consensus (L-Consensus)

and

Registers

We implement L-Consensus through

Obstruction-free Consensus (O-Consensus)

and

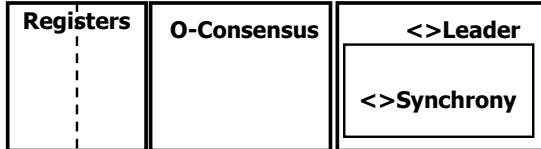
$\langle \rangle$ *Leader* (encapsulating timing assumptions and
sometimes denoted by Ω)

2

A modular approach

Consensus

L-Consensus



3

Consensus

Wait-Free-Termination: If a correct process
proposes, then it eventually decides

Agreement: No two processes decide differently

Validity: Any value decided must have been
proposed

4

L-Consensus

Lock-Free-Termination: If a correct process
proposes, then *at least one* correct process
eventually decides

Agreement: No two processes decide differently

Validity: Any value decided must have been
proposed

5

O-Consensus

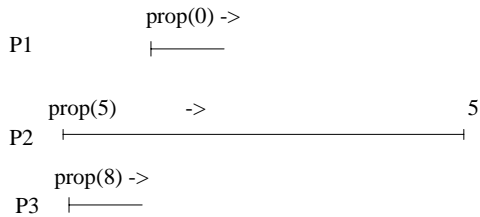
Obstruction-Free-Termination: If a correct process
proposes and *eventually executes alone*, then the
process eventually decides

Agreement: No two processes decide differently

Validity: Any value decided must have been
proposed

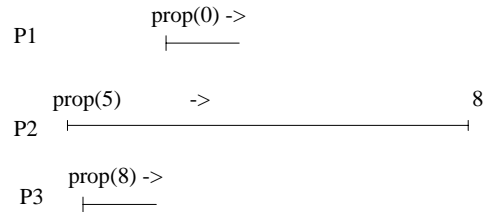
6

Example 1



7

Example 2



8

O-Consensus algorithm (idea)

- ☛ A process that is eventually « left alone » to execute steps, eventually decides
- ☛ Several processes might keep trying to concurrently decide until some unknown time: agreement (and validity) should be ensured during this preliminary period

9

O-Consensus algorithm (data)

- ☛ Each process p_i maintains a timestamp ts , initialized to i and incremented by n
- ☛ The processes share an array of register pairs **Reg[1,..,n]**; each element of the array contains two registers:
 - ☛ **Reg[i].T** contains a timestamp (init to 0)
 - ☛ **Reg[i].V** contains a pair (value,timestamp) (init to (\perp ,0))

10

O-Consensus algorithm (functions)

- ☛ To simplify the presentation, we assume two functions applied to $Reg[1,..,N]$
 - ☛ **highestTsp()** returns the highest timestamp among all elements $Reg[1].T, Reg[2].T, \dots, Reg[N].T$
 - ☛ **highestTspValue()** returns the value with the highest timestamp among all elements $Reg[1].V, Reg[2].V, \dots, Reg[N].V$

11

O-Consensus algorithm

- ☛ `propose(v):`
- ☛ `while(true)`
 - ☛ `Reg[i].T.write(ts);`
 - ☛ `val := Reg[1,..,n].highestTspValue();`
 - ☛ `if val = \perp then val := v;`
 - ☛ `Reg[i].V.write(val,ts);`
 - ☛ `if ts = Reg[1,..,n].highestTsp() then`
 - ☛ `return(val)`
 - ☛ `ts := ts + n`

12

O-Consensus algorithm

- ☞ propose(v):
- ☞ while(true)
 - ☞ (1) Reg[i].T.write(ts);
 - ☞ (2) val := Reg[1,..,n].highestTspValue();
 - ☞ if val = ⊥ then val := v;
 - ☞ (3) Reg[i].V.write(val,ts);
 - ☞ (4) if ts = Reg[1,..,n].highestTsp() then
 - ☞ return(val)
 - ☞ ts := ts + n

13

O-Consensus algorithm

- ☞ (1) pi announces its timestamp
- ☞ (2) pi selects the value with the highest timestamp (or its own if there is none)
- ☞ (3) pi announces the value with its timestamp
- ☞ (4) if pi's timestamp is the highest, then pi decides (i.e., pi knows that any process that executes line 2 will select pi's value)

14

L-Consensus

- ☞ We implement L-Consensus using <>leader (leader()) and the O-Consensus algorithm
- ☞ The idea is to use <>leader to make sure that, eventually, one process keeps executing steps alone, until it decides

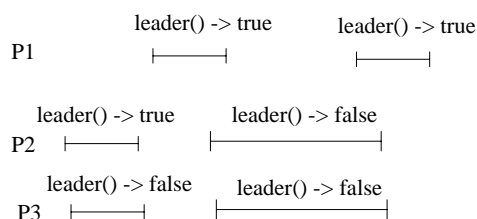
15

<> Leader

- One operation **leader()** which does not take any input parameter and returns, as an output parameter, a boolean
- A process considers itself leader if the boolean is true
- ✓ **Property:** If a correct process invokes leader, then the invocation returns and *eventually*, some correct process is *permanently* the only leader

16

Example



17

L-Consensus

- ☞ propose(v): while(true)
 - ☞ if leader() then
 - ☞ Reg[i].T.write(ts);
 - ☞ val := Reg[1,..,n].highestTspValue();
 - ☞ if val = ⊥ then val := v;
 - ☞ Reg[i].V.write(val,ts);
 - ☞ if ts = Reg[1,..,n].highestTsp()
 - ☞ then return(val)
 - ☞ ts := ts + n

18

From L-Consensus to Consensus (helping)

- Every process that decides writes its value in a register **Dec** (init to \perp)
- Every process periodically seeks for a value in **Dec**

19

Consensus

```

☞ propose(v)
☞ while (Dec.read() =  $\perp$ )
☞ if leader() then
    ☞ Reg[i].T.write(ts);
    ☞ val := Reg[1,..,n].highestTspValue();
    ☞ if val =  $\perp$  then val := p;
    ☞ Reg[i].V.write(val,ts);
    ☞ if ts = Reg[1,..,n].highestTsp()
        ☞ then Dec.write(val)
    ☞ ts := ts + n;
return(Dec.read())
    
```

20

<> Leader

- One operation **leader()** which does not take any input parameter and returns, as an output parameter, a boolean
 - A process considers itself leader if the boolean is true
- ✓ **Property.** If a correct process invokes leader, then the invocation returns and *eventually*, some correct process is *permanently* the only leader

21

<>Leader: algorithm

- We assume that the system is <>synchronous
 - ✓ There is a time after which there is a lower and an upper bound on the delay for a process to execute a local action, a read or a write in shared memory
 - ✓ The time after which the system becomes synchronous is called the global stabilization time (GST) and is unknown to the processes
- This model captures the practical observation that distributed systems are usually synchronous and sometimes asynchronous

22

<>Leader: algorithm (shared variables)

- Every process p_i elects (stores in a local variable leader) the process with the lowest identity that p_i considers as non-crashed; if p_i elects p_j , then $j < i$
- A process p_i that considers itself leader keeps incrementing **Reg[i]**; p_i claims that it wants to remain leader
- NB. Eventually, only the leader keeps incrementing the shared leader

23

<>Leader: algorithm (local variables)

- Every process periodically increments local variables **clock** and **check**, as well as a local variable **delay** whenever its leader changes
- Process p_i maintains **lasti[j]** to record the last value of **Reg[j]** p_i has read (p_i can hence know whether p_j has progressed)
- The next leader is the one with the smallest id that makes some progress; if no such process p_j such that $j < i$ exists, then p_i elects itself (**noLeader** is true)

24

<>Leader: algorithm (variables)

- **check**, and **delay** are initialized to 1
- **lasti[j]** and **Reg[j]** are initialized to 0
- The next leader is the one with the smallest id that makes some progress; if no such process p_j such that $j < i$ exists, then p_i elects itself (**noLeader** is true)

25

<>Leader: algorithm

leader(): return(leader)

- check, delay and leader init to 1
- lasti[j] and Reg[j] init to 0;
- Task:
 - while(true) do
 - ✓ clock := 0;
 - ✓ If (leader=self) then
 - ✓ Reg[i].write(Reg[i].read()+1);
 - ✓ clock := clock + 1;
 - ✓ if(clock = check) then
 - ✓ elect();

26

<>Leader: algorithm (cont'd)

elect():

- noLeader := true;
- for j = 1 to (i-1) do
 - ✓ if (Reg[j].read() > lasti[j]) then
 - ✓ lasti[j] := Reg[j].read();
 - ✓ if(leader ≠ pj) then delay:=delay*2;
 - ✓ check := check + delay;
 - ✓ leader:= pj;
 - ✓ noLeader := false; break (for);
- if (noLeader) then leader := self;

27

Consensus = Registers + <> Leader

- <>Leader has one operation **leader()** which does not take any input parameter and returns, as an output parameter, a boolean; a process considers itself leader if the boolean is true
 - ✓ **Property:** If a correct process invokes leader, then the invocation returns and *eventually*, some correct process is *permanently* the only leader
- <>Leader encapsulates the following synchrony assumption: there is a time after which a lower and an upper bound hold on the time it takes for every process to execute a step (eventual synchrony)

28

Minimal Assumptions

- Consensus is impossible in an asynchronous system with Registers (FLP83, LA88)
- Consensus is possible in an eventually synchronous system (i.e., <> Leader) with Registers (DLS88, LH95)
- What is the minimal synchrony assumption needed to implement Consensus with Registers?
- Is there any weaker timing abstraction than <>Leader that help Registers solve Consensus

29

Failure detector

- A **failure detector** is a distributed (wait-free) oracle that provides processes with information about the **crashes** of processes
- Examples: $P, \diamond P, \diamond S, \diamond W, \Omega, \diamond Leader$
- NB. A failure detector does **only** provide information about crashes (CT96)

30

Failure detector relations

- We say that a failure detector D **implements** abstraction A (e.g., object O) if there is an algorithm that implements A using D
- We say that a failure detector D is **weaker** than a failure detector D' if D' **implements** D ($D \leq D'$)
- If D is weaker than D' and D' is not weaker than D , then D is said to be **strictly weaker** than D' ($D < D'$)
- We say that two failure detectors are **equivalent** if each is weaker than the other ($D \cong D'$)

31

Failure detector Ω

- Failure detector Ω outputs a process q at every process p (we say that p **trusts** q) and ensures the following property:
 - Eventually, the **same correct** process is **permanently trusted** by every process
 - NB. Note that the process that is trusted might keep changing until some eventual time

32

$\langle \rangle$ Leader $\cong \Omega$

- To implement $\langle \rangle$ Leader using Ω , every process simply returns true if it is leader (the process emulates the output of $\langle \rangle$ Leader)
- To implement $\langle \rangle$ Leader using Ω , every process writes its name in a shared register L when leader() returns true; all processes periodically read L and elect the process in L (eventually, only one process is elected)

33

Failure detector example

- Failure detector Ω outputs a process q at every process p (we say that p **trusts** q) and ensures the following property:
 - \diamond **unique leader**: eventually, the **same correct** process is **permanently trusted** by every process
 - NB. Note that the process that is trusted might keep changing until some eventual time

34

Questions

- (1) Show that Ω is the weakest failure detector to implement consensus with Registers (i.e., give an algorithm that implements Ω with any failure detector that implements Consensus with Registers)
- (2) What is the weakest failure detector to implement Consensus with objects of consensus number k and Registers?
- (3) What is the weakest failure to implement an object with consensus number k using Registers?

35