

FDS 2012: Consensus and Set Agreement

© 2012 P. Kuznetsov

So far...

Shared-memory communication:

- safe bits => multi-valued atomic registers
- safe bit => atomic bit (with constant overhead)
- Atomic registers => atomic/immediate snapshot
- Atomic registers \Leftrightarrow Iterated Immediate Snapshot (non-blocking, for tasks)



© 2012 P. Kuznetsov

2

Today

Reaching **agreement** in shared memory:

- Consensus
 - ✓ Impossibility of wait-free consensus
- Set agreement
 - ✓ Impossibility of wait-free consensus (Sperner's lemma)
- 1-resilient consensus impossibility



© 2012 P. Kuznetsov

3

System model

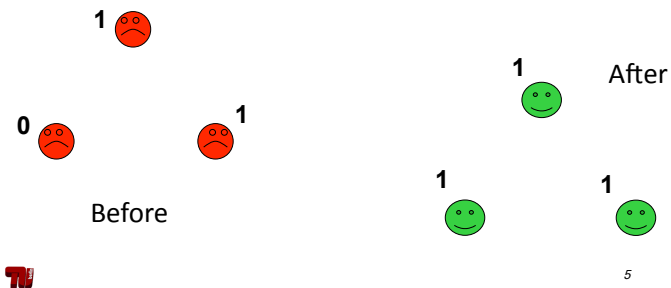
- N *asynchronous* (no bounds on relative speeds) processes p_0, \dots, p_{N-1} ($N \geq 2$) communicate via atomic read-write registers
- Processes can fail by **crashing**
 - ✓ A crashed process takes only finitely many **steps** (reads and writes)
 - ✓ Up to t processes can crash: **t -resilient system**
 - ✓ $t=N-1$: **wait-free**



4

Consensus

Processes *propose* values and must *agree* on a common decision value so that the decided value is a proposed value of some process



Consensus: definition

A process *proposes* an *input* value in V ($|V| \geq 2$) and tries to *decide* on an *output* value in V

- **Agreement:** No two process decide on different values
 - **Validity:** Every decided value is a proposed value
 - **Termination:** No process takes infinitely many steps without deciding
(Every **correct** process decides)
- TU 6

Optimistic (0-resilient) consensus

Consider the case $t=0$, no process fails

Shared: 1WNR register D , initially T (default value not in V)

Upon **propose**(v) by process p_i :

```
if  $i = 0$  then  $D.write(v)$  // if  $p_0$  decide on  $v$   
wait until  $D.read() \neq T$  // wait until  $p_0$  decides  
return  $D$ 
```

(every process decides on p_0 's input)

TU © 2012 P. Kuznetsov 7

Impossibility of wait-free consensus [FLP85,LA87]

Theorem 1 No **wait-free** algorithm solves consensus

We give the proof for $N=2$, assuming that

p_0 proposes 0 and p_1 proposes 1

Implies the claim for all $N \geq 2$

TU 8

Proof of Theorem 1

- We show that no 2-process wait-free solution exists for IIS (iterated immediate snapshot memory)

$r := 0$

repeat

$r := r+1$

$v_i := IS_r.WriteRead_i(v_i)$

until not decided (v_i)

- The IIS=>AS non-blocking simulation (last class) implies that no solution exists for (non-iterated) read-write registers



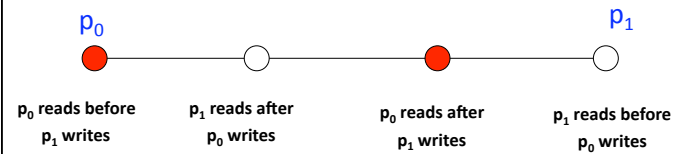
© 2012 P. Kuznetsov

9

Proof of Theorem 1

Initially each p_i only knows its input

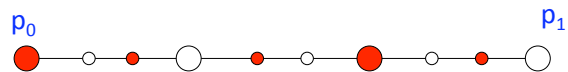
One round of IIS:



10

Proof of Theorem 1

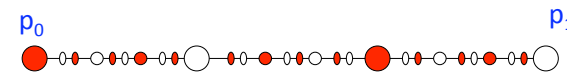
Two rounds of IIS:



11

Proof of Theorem 1

And so on...



Solo runs remain connected - no way to decide!



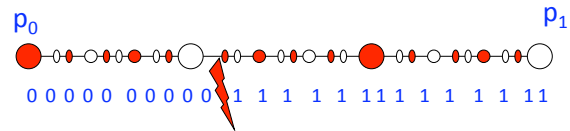
12

Proof of Theorem 1

Suppose p_i ($i=0,1$) proposes i

- p_i must decide i in a solo run!

Suppose by round r every process decides



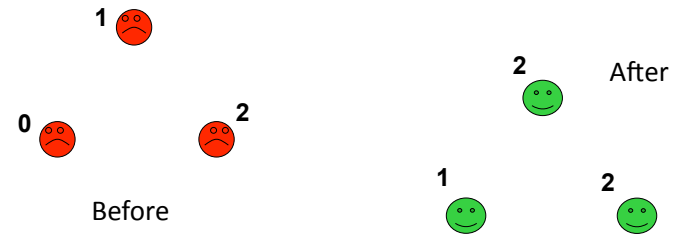
There exists a run with conflicting decisions!



13

Thus

- Wait-free consensus is impossible (for all N)
- What about weaker forms of agreement?
✓ E.g., decide on at most k ($k < N$) different values



© 2012 P. Kuznetsov

14

k-Set Agreement: definition

A process *proposes* an *input* value in V ($|V| \geq k+1$) and tries to *decide* on an *output* value in V

- *k-Agreement*: At most k distinct values are decided
- *Validity*: Every decided value is a proposed value
- *Termination*: No process takes infinitely many steps without deciding
(Every *correct* process decides)

1-set agreement = consensus

($N-1$)-set agreement = **set agreement**



15

Impossibility of Set Agreement

Theorem 2 No **wait-free** algorithm solves set agreement in read-write

- Starting with N values, there is no way to drop one (decide on at most $N-1$)
- Implies the impossibility of wait-free k -set agreement for all $k \leq N-1$



© 2012 P. Kuznetsov

16

Solving Set Agreement in IIS

- Each p_i proposes its **id** i
- There is an IIS round r , such that each process (that reaches round r) outputs **some id**
- In each run, the set of **output ids** is a subset of **participants** of size $\leq N-1$

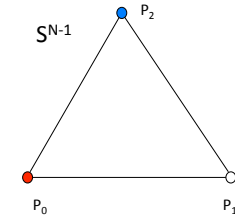
Implies **Sperner coloring** of the subdivided simplex IIS^r



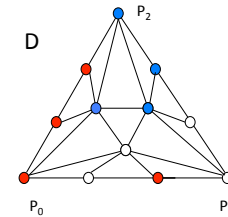
17

Sperner's coloring of a subdivided simplex

Each vertex of the original simplex S^{N-1} (initial state) is **colored** with distinct a process id



Each simplex in **subdivision** D is contained in a face (subset) of S and the union of simplexes of D is S ($|S|=|D|$)



Sperner's coloring:

Each vertex v in the subdivision D is **colored** with colored d :

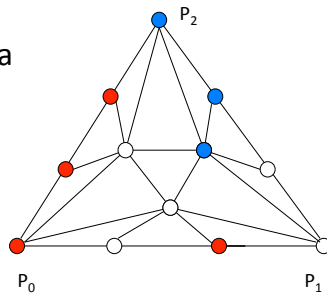
- ✓ If the vertex belongs to a face S of S^{N-1} , then d is in **colors(S)**



18

Sperner's Lemma

Sperner's coloring of any subdivision D of S^{N-1} has a simplex with **all N colors**



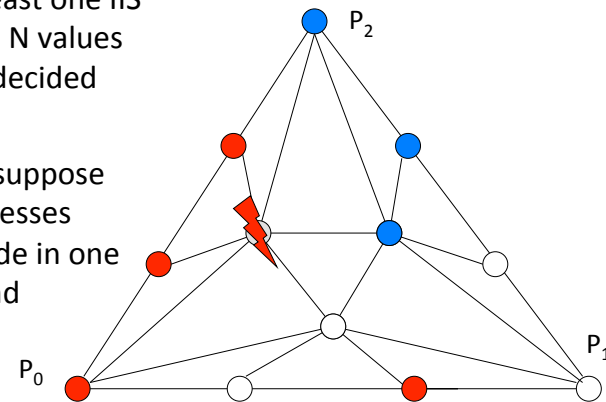
Corollary: wait-free set agreement is impossible in IIS (and, thus, in read-write)



19

In at least one IIS run, N values are decided

Here: suppose processes decide in one round



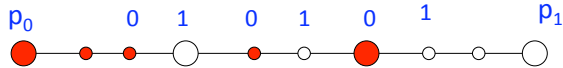
20

Sperner's lemma: inductive step

Claim: for each $k=0, \dots, N-1$, face p_0, \dots, p_{k-1} of S^{N-1} has an **odd** number of k -dimensional simplexes colored $0, \dots, k-1$

By induction: $k=0$ - trivial (exactly one)

$K=1$, simple counting



Suppose the claim holds for $k < N$ and consider the face $1, \dots, k+1$



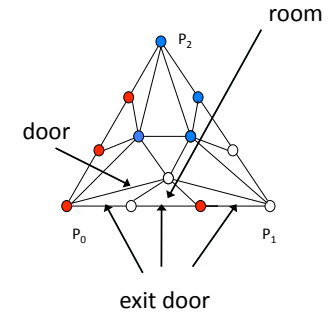
© 2012 P. Kuznetsov

21

Sperner's: rooms and doors

Each k -simplex contained in face p_0, \dots, p_k is a **room**

A $(k-1)$ -dimensional face (a subset of $k-1$ vertices) of a room colored in $0, \dots, k-1$ is a **door**



A door is an **exit** if it is contained in the boundary of p_0, \dots, p_k



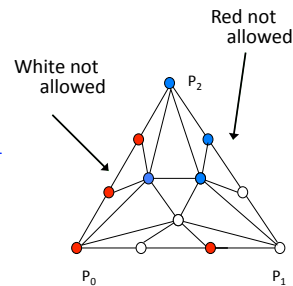
© 2012 P. Kuznetsov

22

Sperner's: rooms and doors

There is an **odd** number of exits!

- By Sperner's coloring no faces other than in p_0, \dots, p_{k-1} can contain simplexes colored $0, \dots, k-1$
- Exits may only be contained in p_0, \dots, p_{k-1}
- By induction, p_0, \dots, p_{k-1} contains an odd number of doors (colored with $0, \dots, k-1$)



© 2012 P. Kuznetsov

23

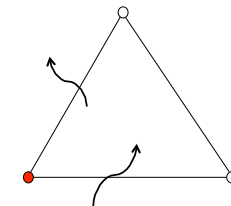
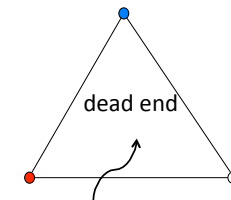
Sperner's: corridors and dead ends

Consider a room with a door (k vertices colored $0, \dots, k-1$)

Two cases possible (depending on the color of the remaining vertex):

- the remaining vertex is colored k - the room has exactly one door (dead end)
 - ✓ the room is fully colored $(0, \dots, k)$
- The remaining vertex is colored in one of $\{0, \dots, k-1\}$ - the room has two doors

We show that there is an odd number of dead ends contained in face p_0, \dots, p_k



© 2012 P. Kuznetsov

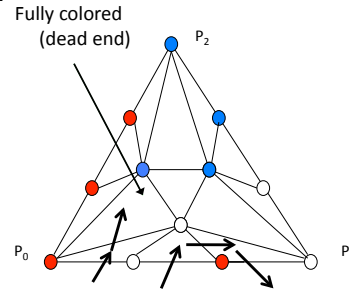
24

Sperner's: counting fully colored rooms

Start with an exit and walk through the doors

There are only finitely many rooms, thus, two only cases possible:

- Stop in a dead end (fully colored simplex)
- Reach another exit



© 2012 P. Kuznetsov

25

Sperner's: crossing the rooms

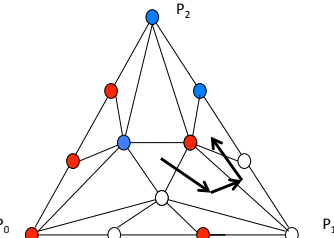
The number of exits is odd

=> The number of dead ends (fully colored simplexes) reachable from exit is **odd!**

Now start in a dead end **not reachable** from an exit: can only stop in another dead end not reachable from an exit

=> The number of dead ends not reachable from exit is **even!**

=> **The total number of fully colored rooms is odd**



© 2012 P. Kuznetsov

26

Thus

- No algorithm can wait-free solve (N-1)-set agreement in IIS
 - ✓ otherwise Sperner's coloring of some IIS^r would have no fully colored simplexes (but there is at least one!)
- No algorithm can wait-free solve set agreement in AS (or any read-write model):
 - ✓ otherwise we can (non-blocking) simulate it in IIS and thus find a wait-free algorithm in IIS
- We cannot tolerate N-1 failures: can we tolerate less?
 - ✓ E.g., can we solve consensus (1-set agreement) **1-resiliently?**



© 2012 P. Kuznetsov

27

1-resilient consensus?

What if we have 1000000 processes and one of them can crash?

NO

We present a direct proof now (an indirect proof by reduction to the wait-free impossibility also exists)



© 2012 P. Kuznetsov

28

Impossibility of 1-resilient consensus [FLP85,LA87]

Theorem 2 No 1-resilient (assuming that one process might fail) algorithm solves consensus in read-write

Proof

By contradiction, suppose that an algorithm A solves 1-resilient binary consensus among p_0, \dots, p_{N-1}



29

Proof

By contradiction, suppose that an algorithm A solves wait-free binary consensus among p_0, \dots, p_{N-1}

A run of A is a sequence of atomic steps (reads or writes) applied to the initial state

A run of A can be seen as an initial input configuration (one input per process) and a sequence of process ids $i_1, i_2, \dots, i_k, \dots$ (all registers are atomic)

Every correct (taking sufficiently many steps) process decides!



30

Proof: valence

Let R be a finite run

- We say that R is *v-valent* (for v in $\{0,1\}$) if v is decided in every infinite extension of R
- We say that R is *bivalent* if R is neither 0-valent nor 1-valent (there exists a 0-valent extension of R and a 1-valent extension of R)



31

Proof: valence claims

Claim 1 Every finite run is 0-valent, or 1-valent, or bivalent. (by Termination)

Claim 2 Any run in which some process decides v is v -valent (by Agreement)

Corollary 1: No process can decide in a bivalent run (by Agreement).



© 2012 P. Kouznetsov

32

Bivalent input

Claim 3 There exists a bivalent input configuration (empty run)

Proof

Suppose not

Consider sequence of input configurations C_0, \dots, C_N :

C_i : p_0, \dots, p_{i-1} propose 1, and p_i, \dots, p_{N-1} propose 0

- All C_i 's are univalent
- C_0 is 0-valent (by Validity)
- C_N is 1-valent (by Validity)



© 2012 P. Kuznetsov

33

Bivalent input

There exists i in $\{0, \dots, N-2\}$ such that C_i is 0-valent and C_{i+1} is 1-valent!

C_i and C_{i+1} differ **only in the input value of p_i** (it proposes 1 in C_i and 0 in C_{i+1})

Consider a run R starting from C_i in which p_i takes no steps (crashes initially): eventually all other processes decide 0

Consider R' that is like R except that it starts from C_{i+1}

- R and R' are **indistinguishable!**
- Thus, every process decides 0 in R' --- contradiction (C_{i+1} is 1-valent)



© 2012 P. Kuznetsov

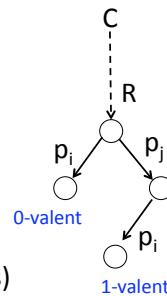
34

Critical run

Claim 4 There exists a **critical** (bivalent) finite run R and two processes p_i and p_j such that **$R.i$ is 0-valent** and **$R.j.i$ is 1-valent** (or vice versa)

Proof of Claim 4: By construction, take the bivalent empty run C (by Claim 3 it exists)

We construct an ever-extending fair (giving each process enough steps) run which results in R



© 2012 P. Kuznetsov

35

Proof of Claim 4: critical run

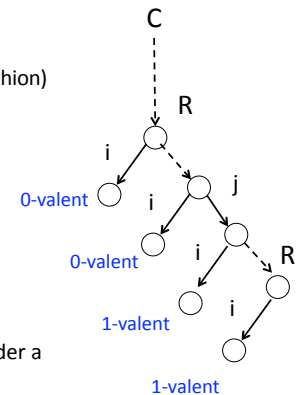
repeat forever

take the **next** process p_i (in **round-robin** fashion)

if for some R' , an extension of R , $R.i$ is bivalent **then** $R := R'.i$

else stop

- If never stops – ever extending (infinite) bivalent runs in which every process is correct (takes infinitely many steps – contradiction with termination)
- If stops – (suppose $R.i$ is 0-valent) – consider a 1-valent extension
 - ✓ There is a critical configuration between R and R'



© 2012 P. Kuznetsov

36

Proof (contd.)

Take a critical run R (exists by Claim 4) such that:

- $R.0$ is 0-valent
- $R.1.0$ is 1-valent

(without loss of generality, we can always rename processes or inputs appropriately ☺)



© 2012 P. Kouznetsov

37

Proof (contd.): the next steps in R

Four cases, depending on the next steps of p_0 and p_1 in R

- p_0 and p_1 are about to access different objects in R
- p_1 reads X and p_0 reads X
- p_0 writes in X
- p_1 reads X

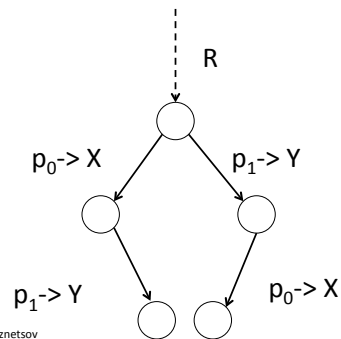


© 2012 P. Kouznetsov

38

Proof (contd.): cases and contradiction

- p_0 and p_1 are about to access **different** objects in R
 ✓ $R.0.1$ and $R.1.0$ are indistinguishable



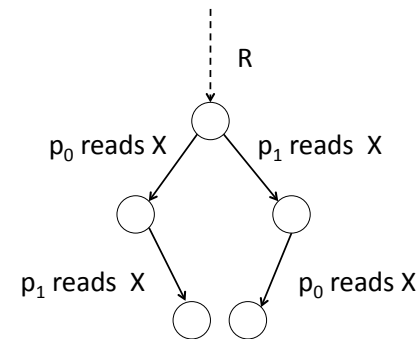
© 2012 P. Kouznetsov



39

Proof (contd.): cases and contradiction

- p_0 and p_1 are about to **read** the same object X
 $R.0.1$ and $R.1.0$ are indistinguishable

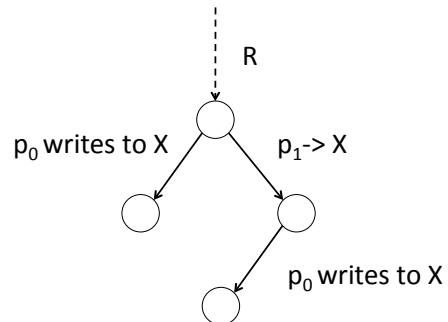


© 2012 P. Kouznetsov

40

Proof (contd.): cases and contradiction

- p_0 is about to write to X
 - ✓ Extensions of R.0 and R.1.0 are indistinguishable for all except p_1 (assuming p_1 takes no more steps)

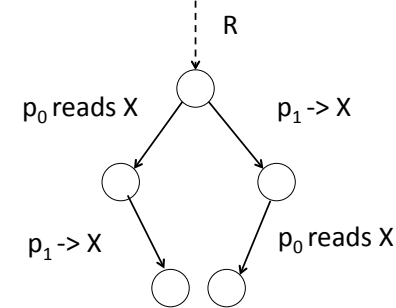


© 2012 P. Kouznetsov

41

Proof (contd.): cases and contradiction

- p_0 is about to read to X
 - ✓ Extensions of R.0.1 and R.1.0 are indistinguishable for all but p_0 (assuming p_0 takes no more steps)



© 2012 P. Kouznetsov

42

So...

- No critical run exists
 - A contradiction with **Claim 4**
- ⇒ 1-resilient consensus is impossible in read-write



© 2012 P. Kouznetsov

43

- 1-resilient consensus is impossible in read-write
- Can be generalized to:
 - ✓ k-resilient k-set agreement is impossible in read-write
 - ✓ (by a reduction to the wait-free case)

Next class:

- showing that consensus is **universal** (wait-free implements any sequential object)
- Solving consensus with stronger (than read-write) objects

Next tutorial: Q&A (homework: go through the previous homework 😊)



© 2012 P. Kouznetsov

44