

institute **imdea**

institute **imdea**
networks

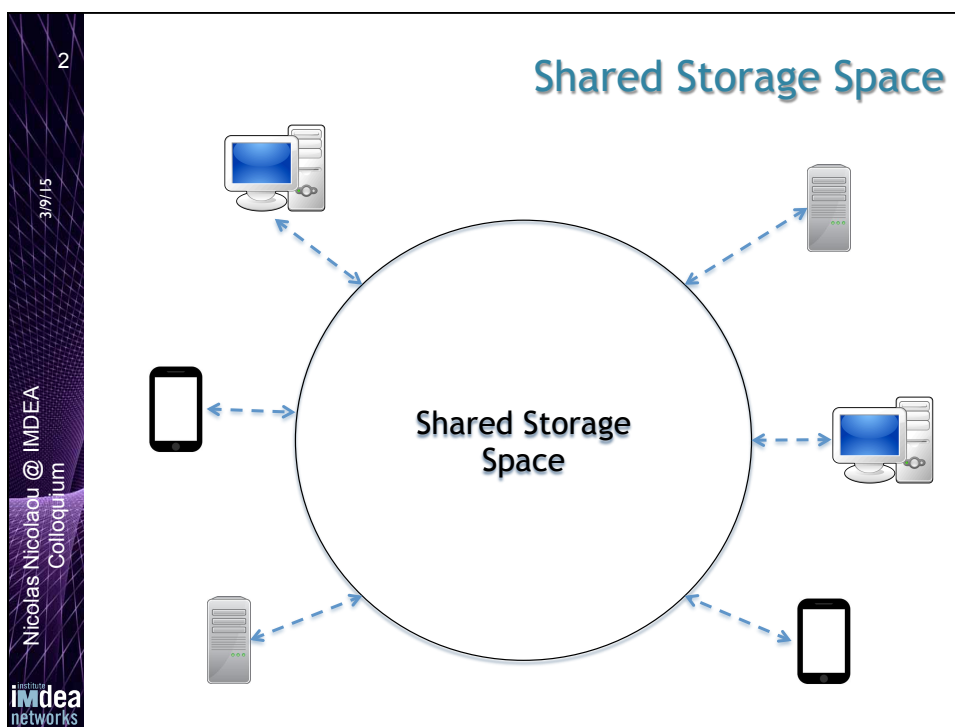





Emulating Highly Consistent Shared Storage on top of Unreliable Message-Passing Nodes

Nicolas Nicolaou

[Developing the
Science of Networks]



3

3/9/15

Nicolas Nicolaou @ IMDEA Colloquium

imdea networks

Distributed Storage Solution

Alice $\text{write}(y)$

Bob $\text{read}()$

Distributed Storage Abstraction

- Data Replication – Servers/Disks
 - Survivability and Availability
- Read/Write operations
- Consistency Semantics

4

3/9/15

Nicolas Nicolaou @ IMDEA Colloquium

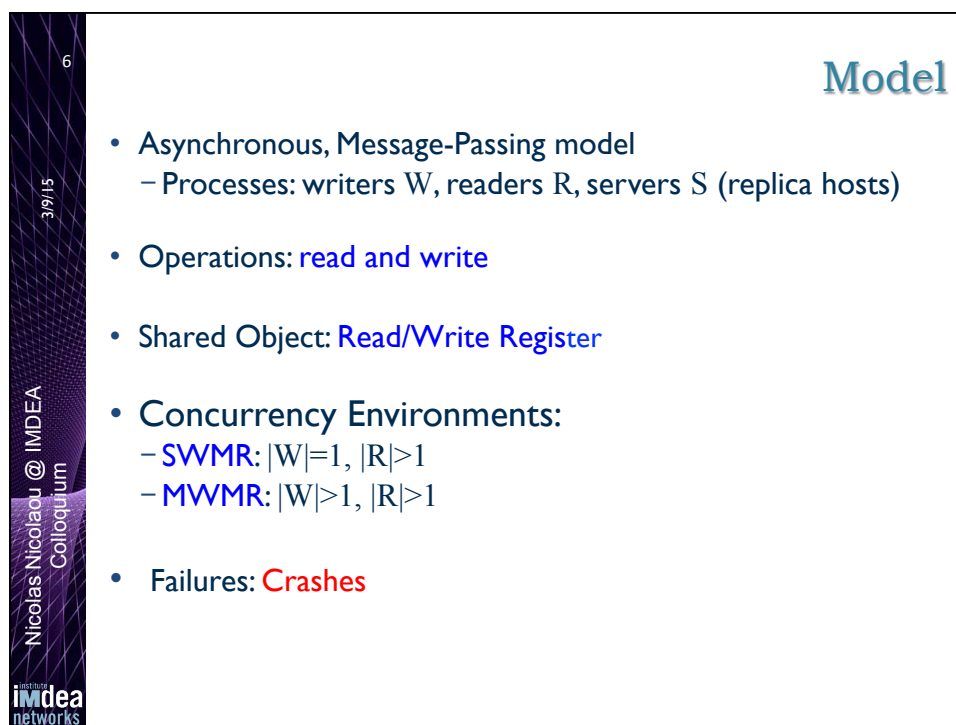
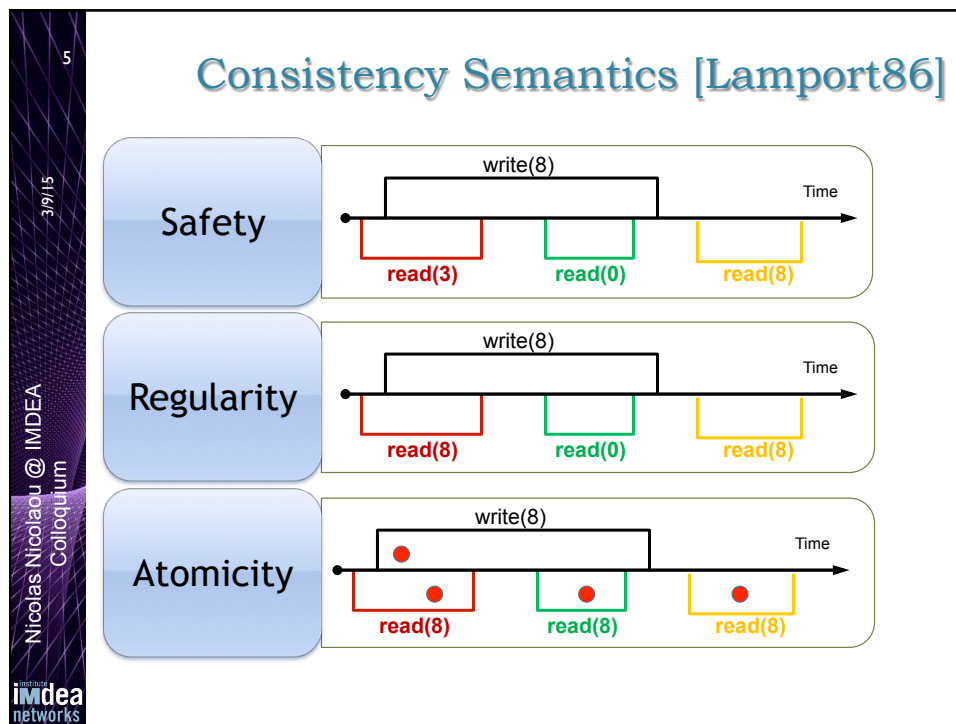
imdea networks

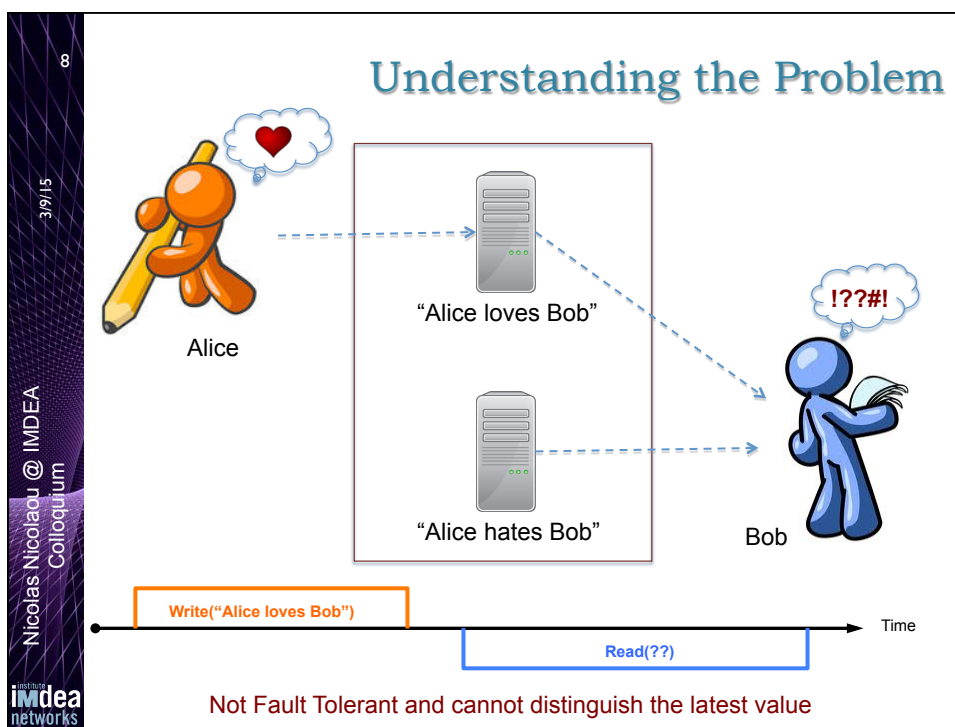
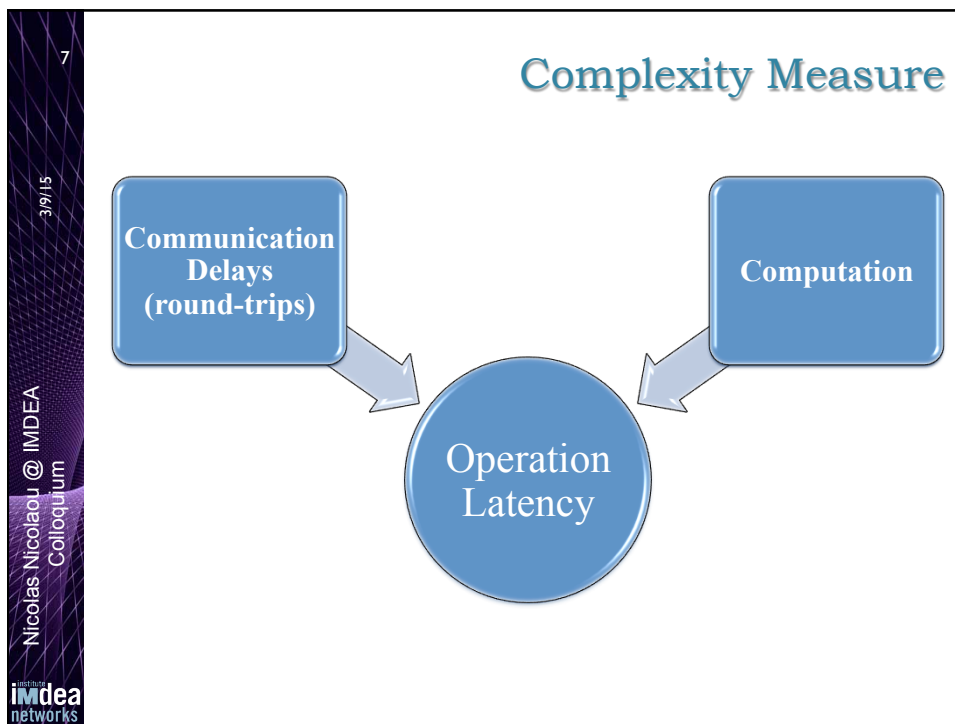
Operation Relations

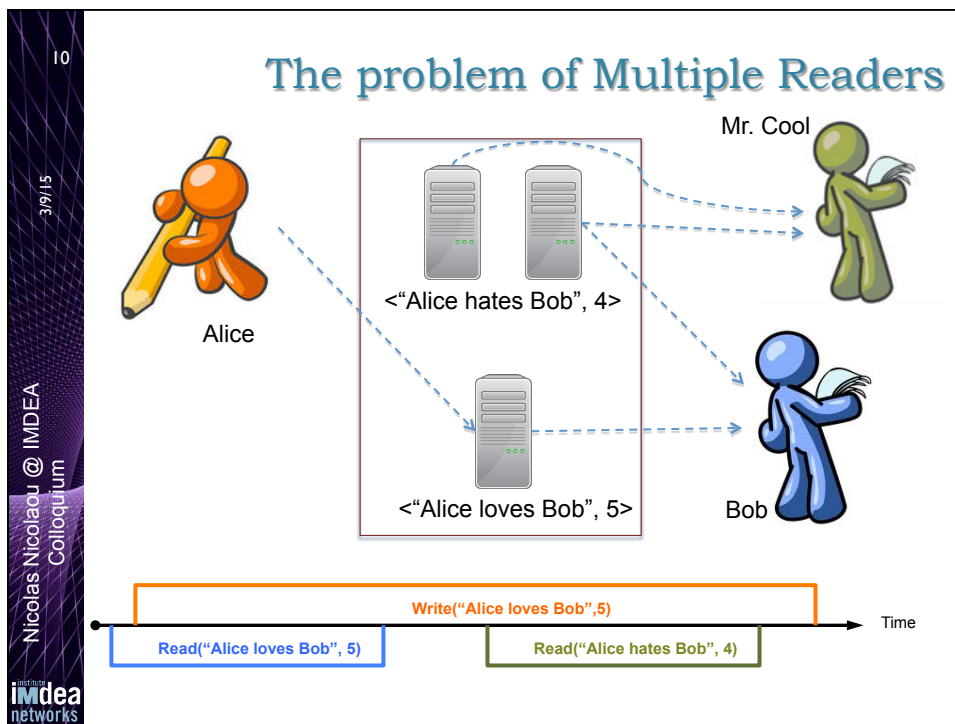
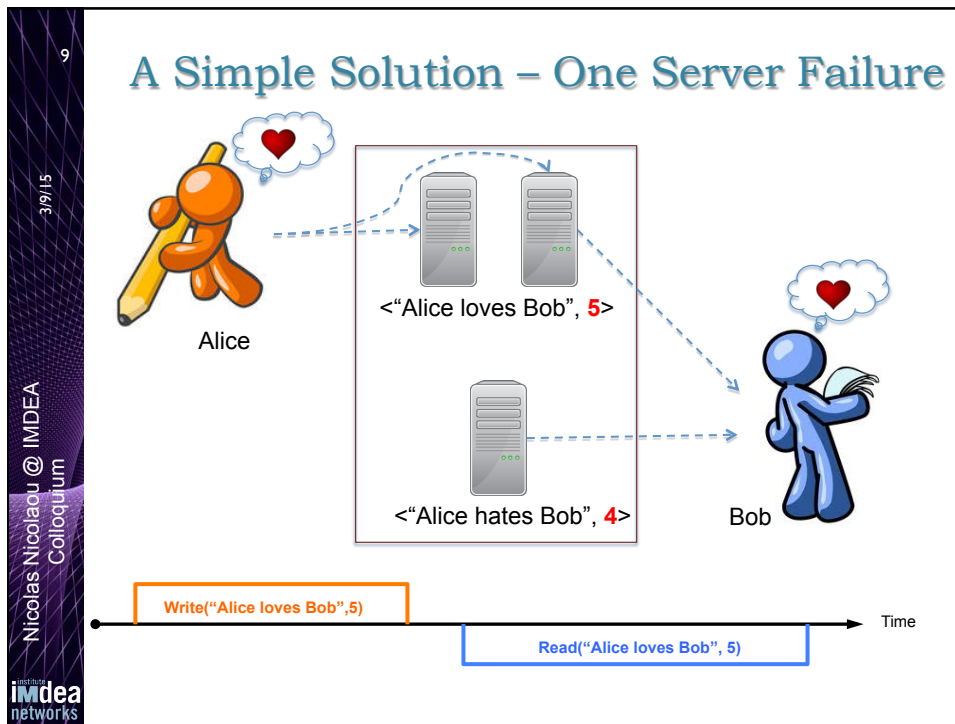
- Precedence Relations for two operations π_1, π_2 :
 - π_1 **precedes** π_2 if the response of π_1 happens **before** the invocation of π_2

- π_1 **succeeds** π_2 if the invocation of π_1 happens **after** the response of π_2

- π_1 is **concurrent** with π_2 if π_1 **neither precedes nor succeeds** π_2







[Attiya, Bar-Noy, Dolev '96]

Algorithm: ABD

Write Protocol: one phase

- P1: Increment your local timestamp, and send $\langle ts, v \rangle$ to a majority of the replicas

Read Protocol: two phases

- P1: **query phase**
 - Send $\text{read}()$ to all the replicas
 - When a majority replies discover the largest $\langle \text{maxts}, v \rangle$
- P2: **propagation phase**
 - Send $\text{write}(\langle \text{maxts}, v \rangle)$ to all replicas and wait for a majority to reply

Server Protocol: passive role

- Receive requests, update local timestamp (if $\text{msg.ts} > \text{server.ts}$) and reply with $\langle \text{server.tag}, v \rangle$

11

3/9/15

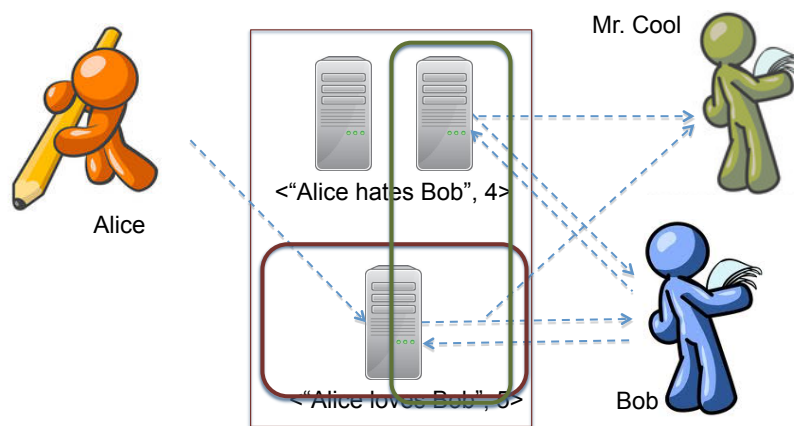
Nicolas Nicolaou @ IMDEA
Colloquiumimdea
networks

12

3/9/15

Nicolas Nicolaou @ IMDEA
Colloquiumimdea
networks

How it works in our example?



Folklore belief: "Reads must Write"

13

3/9/15

Nicolas Nicolaou @ IMDEA Colloquium

imdea networks

The Era of Fast Implementations....

SWMR Fast
[Dutta et al. 2004]

- Single round (**fast**) writes and reads
- **Bounded readers:** $R < (S/f) - 2$ where S servers & f failures
- **Impossible in MWMR model**

SWMR Semifast
[Georgiou, Nicolaou, Shvartsman 2006]

- Single round writes
- Only a **single** complete 2-round (**slow**) read per write
- **Unbounded readers**
- **Bounded Virtual Nodes:** $V < (S/f) - 2$
- **Impossible in the MWMR model**

SWMR Weak-Semifast
[Georgiou, Nicolaou, Shvartsman 2008]

- **General Quorum System**
- Fast writes and Multiple slow reads per write
- Allows concurrent fast reads with writes
- **Unknown if applicable in MWMR model**

14

3/9/15

Nicolas Nicolaou @ IMDEA Colloquium

imdea networks

What happens in the MWMR?

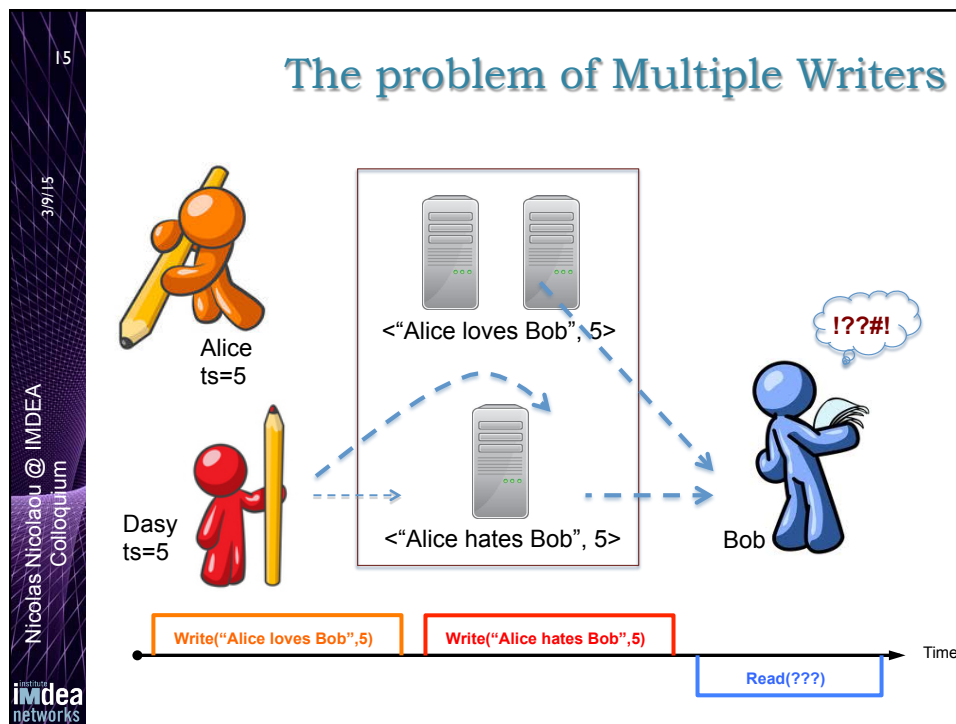
Alice
ts=5

Dasy
ts=5

Bob

<"Alice loves Bob", 5>

<"Who loves Bob?", 4>



16

3/9/15

Nicolas Nicolaou @ IMDEA Colloquium

imdea networks

How to solve the problem?

- Writer needs to discover the **latest timestamp** in the system before incrementing it
- Each writer includes its id with the timestamp
 - Breaks the symmetry in case two writers read the same maximum timestamp
 - The <timestamp, id> pair is called **TAG**
 - tag1 > tag2 if either:
 - tag1.timestamp > tag2.timestamp, or
 - tag1.timestamp = tag2.timestamp AND tag1.id > tag2.id

17

[Lynch, Shvartsman 1997]

Algorithm: MWABD

Belief: "Writes must Read"

Write Protocol: two phases

- P1: **Query** the majority for the largest tag
- P2: **Increment the largest tag**, and **propagate** $\langle \text{newtag}, v \rangle$ to a majority

Read Protocol: two phases

- P1: **Query** the majority for the largest $\langle \text{maxtag}, v \rangle$
- P2: **Propagate** $\langle \text{maxtag}, v \rangle$ to a majority

Server Protocol: passive role

- Receive requests, update local timestamp (if $\text{msg.ts} > \text{server.ts}$) and reply with $\langle \text{server.tag}, v \rangle$

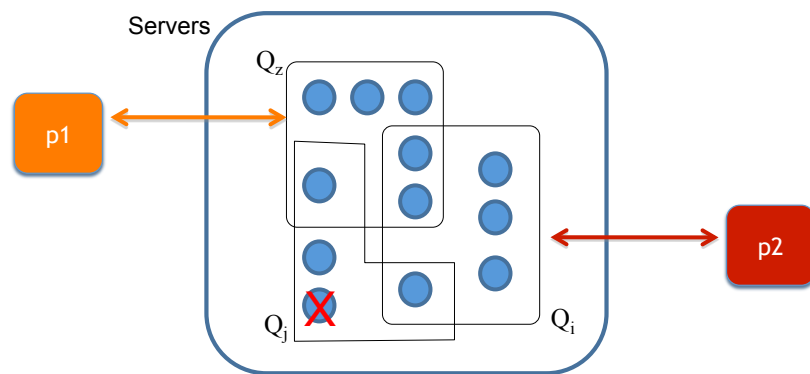
3/9/15

Nicolas Nicolaou @ IMDEA Colloquium

imdea
networks

18

Definition: Quorum systems



- Q_p, Q_i, Q_z are **quorums**
- **Quorum System** is the set $\{Q_p, Q_i, Q_z\}$
 - Property: every pair of quorums intersects
 - **N-wise quorums systems**: every N quorums intersect for $N > 1$
- Every R/W operation communicates with a single quorum
- **Faulty Quorum**: Contains a faulty process

3/9/15

Nicolas Nicolaou @ IMDEA Colloquium

imdea
networks

19

3/9/15

Nicolas Nicolaou @ IMDEA Colloquium

imdea networks

[Englert, Georgiou, Musial, Nicolaou, Shvartsman 2009]

Are Fast Operations Possible in MWMR?

Theorem: No execution of safe register implementation that use an N -wise quorum system, contains more than $N - 1$ consecutive, quorum switching, fast writes.

Theorem: It is impossible to obtain MWMR safe register implementations that exploit an N -wise quorum system, if

$$|W \cup R| > N - 1$$

20

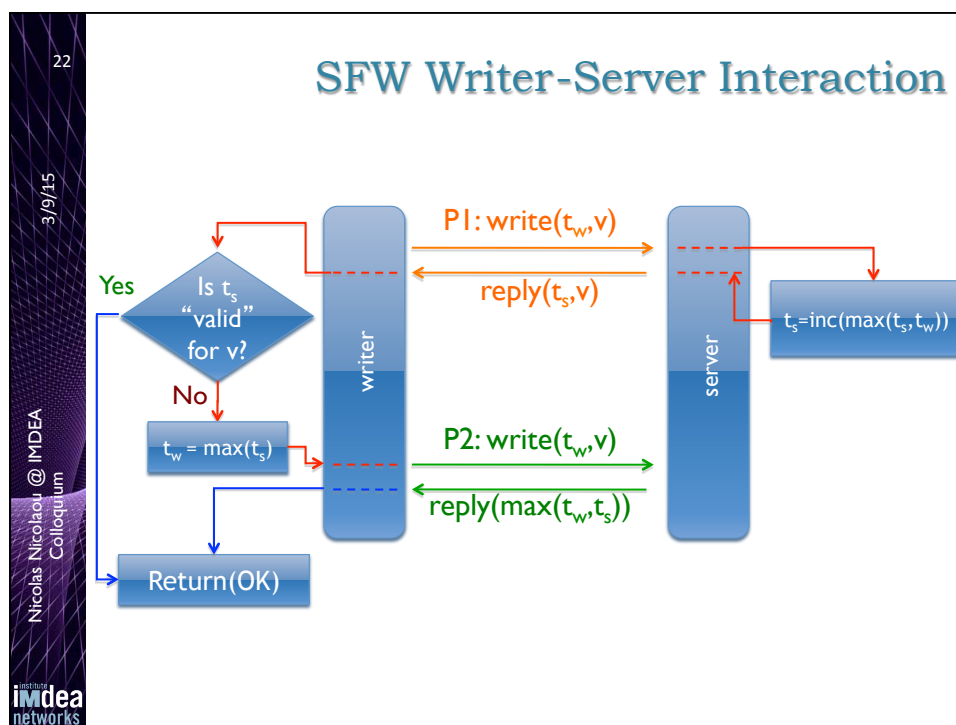
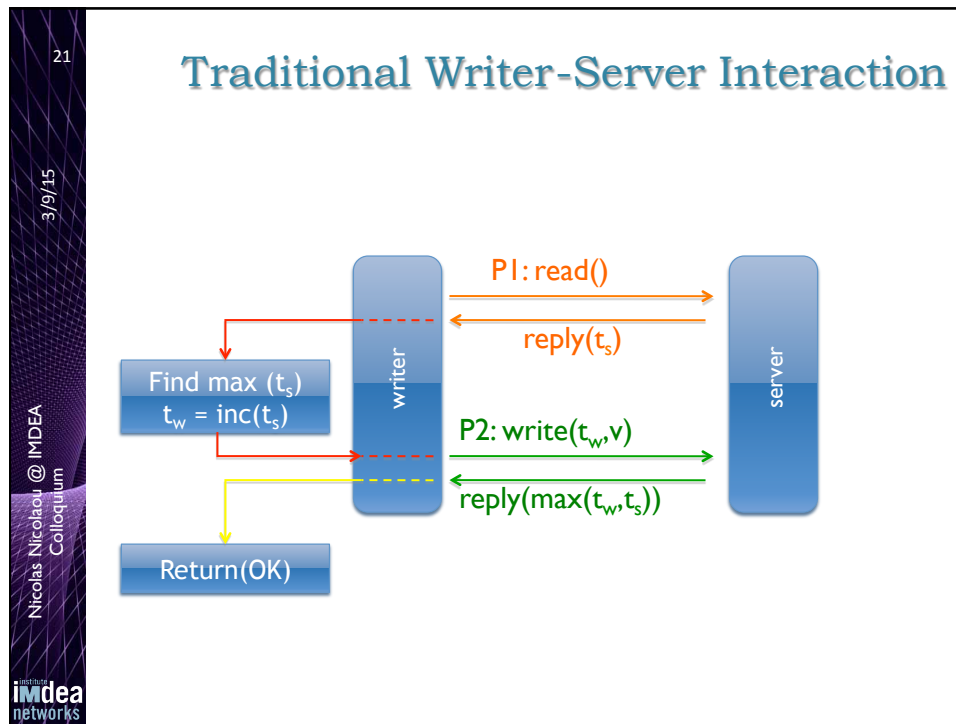
3/9/15

Nicolas Nicolaou @ IMDEA Colloquium

imdea networks

New Technique - SSO

- SSO: Server Side Ordering
 - Tag is incremented by the servers and not by the writer.
 - Generated tags may be different across servers
 - Clients decide operation ordering based on server responses
- SFW Algorithm
 - Enables FastWrites and Reads -- first such algorithm
 - Allows Unbounded Participation



23

[Englert, Georgiou, Musial, Nicolaou, Shvartsman 2009]

Algorithm: SFW (in a glance)

Write Protocol: one or two rounds

- P1: Collect **candidate** tags from a quorum
- Exists tag t propagated in a **bigger** than $(n/2-1)$ -wise intersection (PREDICATE PW)
 - **YES** – assign t to the written value and return => **FAST**
 - **NO** – propagate the unique **largest tag** to a quorum => **SLOW**

Read Protocol: one or two rounds

- P1: collect **list of writes** and their tags from a quorum
- Exists max write tag t in a **bigger** than $(n/2-2)$ -wise intersection (PREDICATE PR)
 - **YES** – return the value written by that write => **FAST**
 - **NO** – is there a confirmed tag propagated to $(n-1)$ -wise intersection => **FAST**
 - **NO** – propagate the largest confirmed tag to a quorum => **SLOW**

Server Protocol

- **Increment tag** when receive write request and send to read/write **the latest writes**

24

Predicates: Read and Write

Writer predicate for a write ω (PW): $\exists \tau, Q^i, MS$ where: $\tau \in \{\langle \cdot, \omega \rangle : \langle \cdot, \omega \rangle \in m(\omega)_{s,w}.inprogress \wedge s \in Q\}$, $MS = \{s : s \in Q \wedge \tau \in m(\omega)_{s,w}.inprogress\}$, and $Q^i \subseteq Q, 0 \leq i \leq \lfloor \frac{n}{2} - 1 \rfloor$, s.t. $(\bigcap_{Q \in Q^i \cup \{Q\}} Q) \subseteq MS$.

Reader predicate for a read ρ (PR): $\exists \tau, Q^j, MS$, where: $\max(\tau) \in \bigcup_{s \in Q} m(\rho)_{s,r}.inprogress$, $MS = \{s : s \in Q \wedge \tau \in m(\rho)_{s,r}.inprogress\}$, and $Q^j \subseteq Q, 0 \leq j \leq \lfloor \frac{n}{2} - 2 \rfloor$, s.t. $(\bigcap_{Q \in Q^j \cup \{Q\}} Q) \subseteq MS$.

25

The Weak Side of SFW

- Predicates are Computationally Hard
 - NP-Complete
- Restriction on the Quorum System
 - Deploys N-wise Quorum Systems
 - Guarantees fastness iff $n > 3$

Nicolas Nicolaou @ IMDEA Colloquium

imdea networks

26

The Good News...

- Algorithm CWFR
 - Based on Quorum Views
 - SWMR prediction tools
 - Fast operations in General Quorum Systems
 - Trades Speed of Write operations
 - Two Round Writes

Nicolas Nicolaou @ IMDEA Colloquium

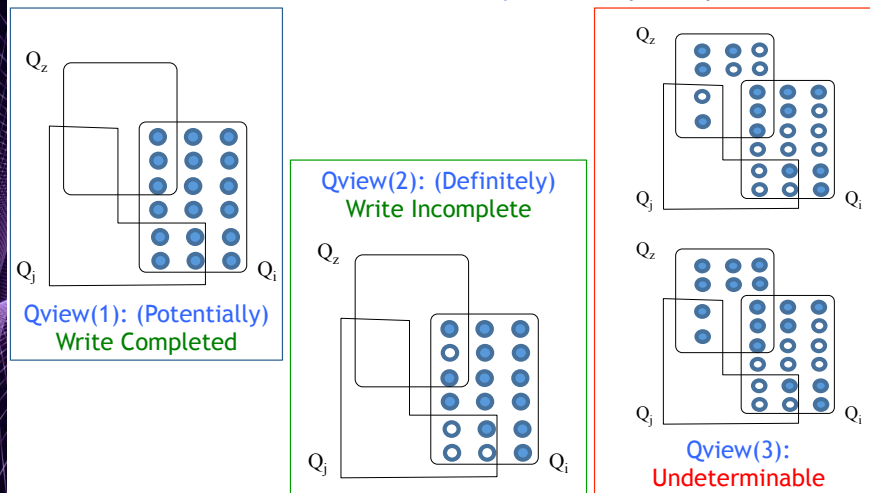
imdea networks

27

[Georgiou, Nicolaou, Shvartsman 2008]

Quorum Views

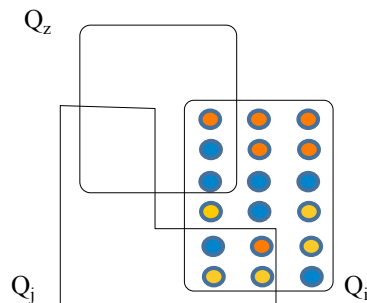
- Idea: Try to **determine the status of the write** operation based on the **distribution of the max timestamp** in the replied quorum.



28

What happens in MWMR?

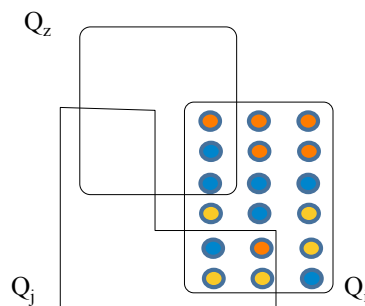
- MWMR environment
 - Concurrent writes
 - Multiple **concurrent values**
- For values $\langle \text{tag1}, v1 \rangle$, $\langle \text{tag2}, v2 \rangle$, $\langle \text{tag3}, v3 \rangle$
 - Let $\text{tag1} < \text{tag2} < \text{tag3}$



29

Idea: Uncover the Past

- Discover the **latest potentially completed** write
- For values $\langle \text{tag1}, v1 \rangle$, $\langle \text{tag2}, v2 \rangle$, $\langle \text{tag3}, v3 \rangle$:
 - $\langle \text{tag3}, v3 \rangle$ **not completed** (servers **possibly** contained $\langle \text{tag2}, v2 \rangle$)
 - $\langle \text{tag2}, v2 \rangle$ **not completed** (servers **possibly** contained $\langle \text{tag1}, v1 \rangle$)
 - $\langle \text{tag1}, v1 \rangle$ **potentially completed**



30

[Georgiou, Nicolaou, Russell, Shvartsman 2012]

Algorithm: CWFR

Traditional Write Protocol: two rounds

- P1: Query a single quorum for the latest tag
- P2: Increment the max tag, send $\langle \text{newtag}, v \rangle$ quorum

Read Protocol: one or two rounds

- Iterate to discover smallest completed write
- P1: receive replies from a quorum Q
 - $QView_Q(1)$ - **Fast**: return maxTag of current iteration
 - $QView_Q(2)$ - **remove servers with maxTag and re-evaluate**
 - $QView_Q(3)$ - **Slow**: propagate and return maxTag_0

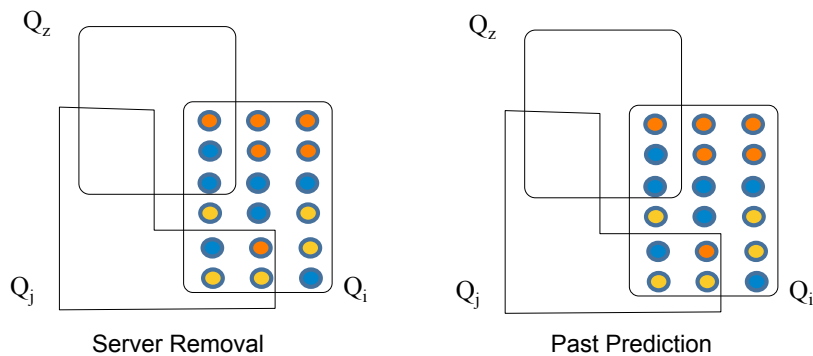
Server Protocol: passive role

- Receive requests, update local timestamp and return $\langle \text{tag}, v \rangle$

31

Read Iteration: Discard Incomplete Tags

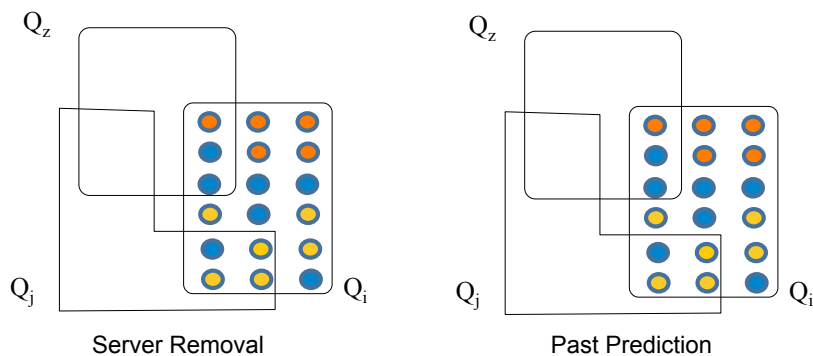
- For values $\langle \text{tag1}, v1 \rangle$, $\langle \text{tag2}, v2 \rangle$, $\langle \text{tag3}, v3 \rangle$:
 - $\langle \text{tag3}, v3 \rangle$ not completed: remove servers that contain $\langle \text{tag3}, v3 \rangle$
 - $\langle \text{tag2}, v2 \rangle$ not completed: remove servers that contain $\langle \text{tag2}, v2 \rangle$
 - $\langle \text{tag1}, v1 \rangle$ potentially completed in Q_i
 - $Q_{\text{view}}(1)$: all remaining servers contain $\langle \text{tag1}, v1 \rangle$



32

Read Iteration: Discard Incomplete Tags

- For values $\langle \text{tag1}, v1 \rangle$, $\langle \text{tag2}, v2 \rangle$, $\langle \text{tag3}, v3 \rangle$:
 - $\langle \text{tag3}, v3 \rangle$ not completed: remove servers that contain $\langle \text{tag3}, v3 \rangle$
 - $\langle \text{tag2}, v2 \rangle$ potentially completed in Q_j
 - $Q_{\text{view}}(3)$: an intersection of the remaining servers contains $\langle \text{tag2}, v2 \rangle$
 - P2: propagate $\langle \text{tag3}, v3 \rangle$ to a complete quorum (help $\langle \text{tag3}, v3 \rangle$ to complete)



33

3/9/15

Nicolas Nicolaou @ IMDEA
Colloquiumimdea
networks

Empirical Evaluation

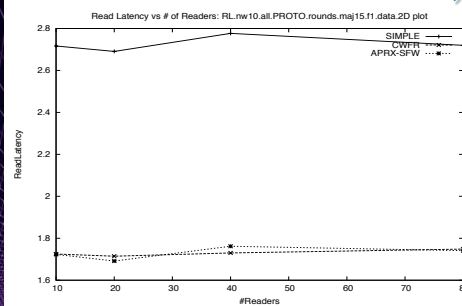
- NS2: Single Processor Emulator
 - Fully Controlled Environment
 - Tested the efficiency of the algorithms under predefined environmental parameters
- Planet-Lab: Planetary Scale Real Time Platform
 - Real Time Network Environment (INTERNET)
 - Machines Disseminated Throughout the Globe
 - Tested practicality of the algorithms

34

3/9/15

Nicolas Nicolaou @ IMDEA
Colloquiumimdea
networks

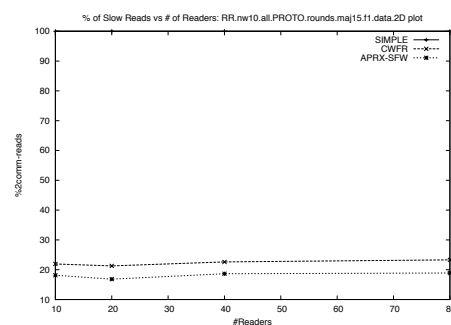
NS2: SIMPLE, APRX-SFW, CWFR

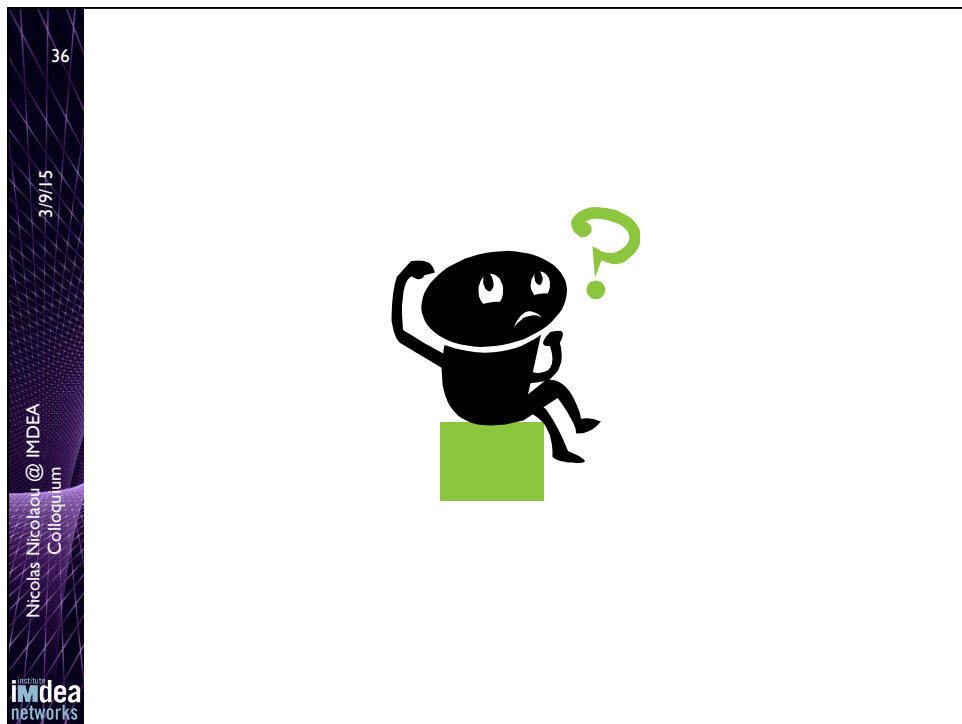
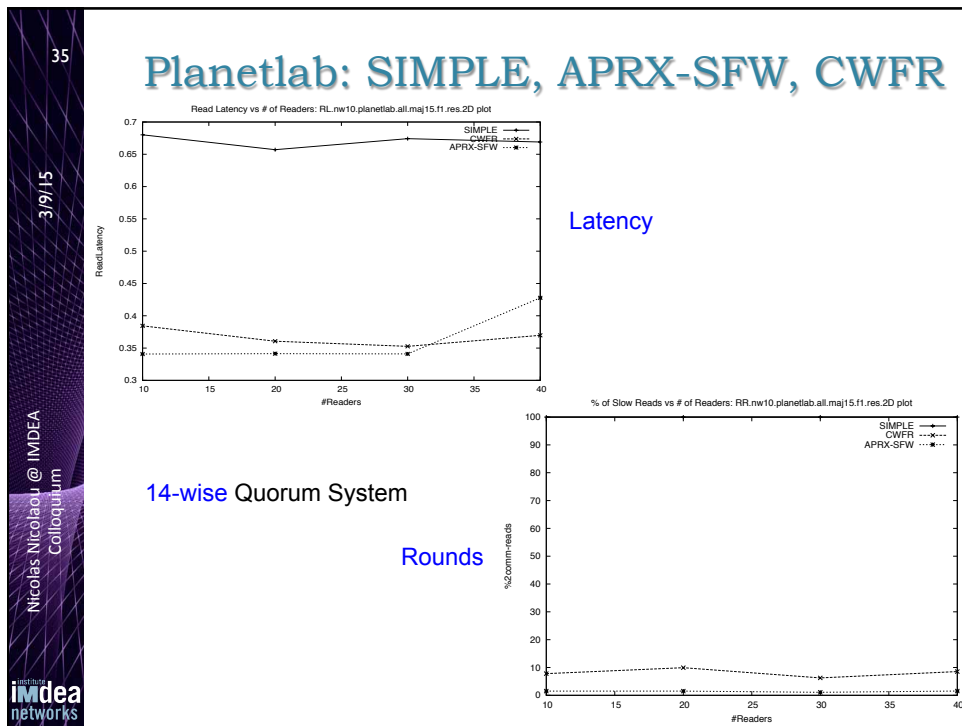


Latency

14-wise Quorum System

Rounds





37

Survivability of Data

• Data Survivability is **essential** in today's systems and applications

• Popular approach: Redundancy

– Use of a Redundant Array of Inexpensive Disks (RAID)

– A RAID system contains

- A single box
- Residing on a single physical location
- Servers are connected to clients via a single network interface

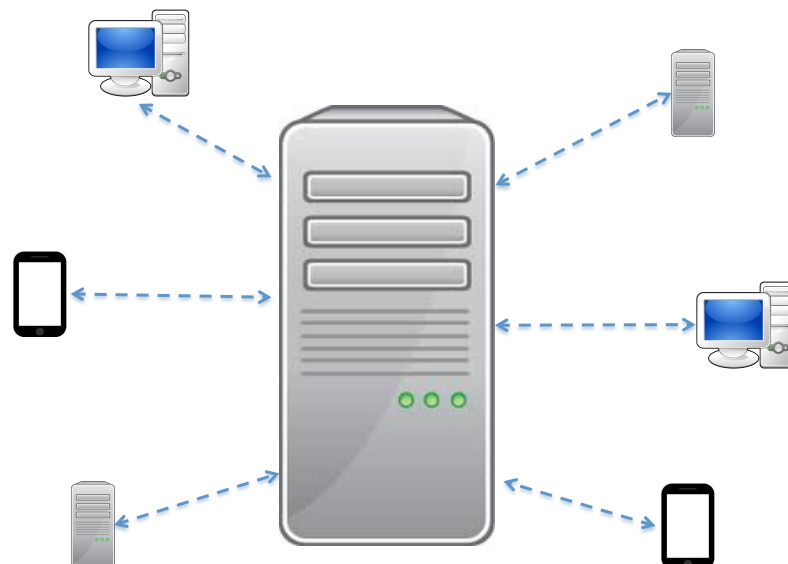
Single point of failure

Use Distributed Storage Systems instead

imdea
networks

38

Centralized Solution??

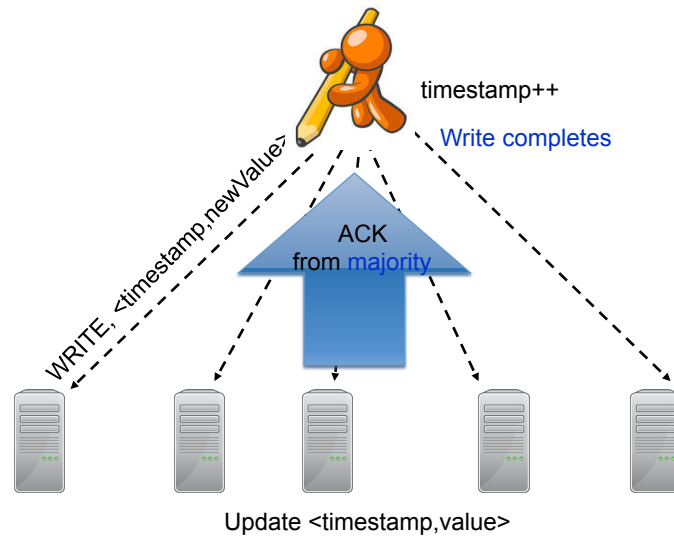


Nicolas Nicolaou @ IMDEA
Colloquium

imdea
networks

39

A Simple Algorithm – ABD: Writer

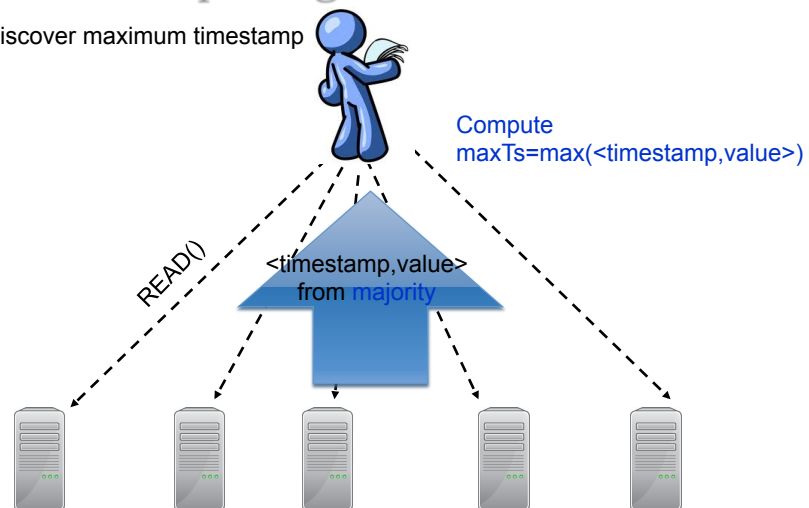


imdea
networks

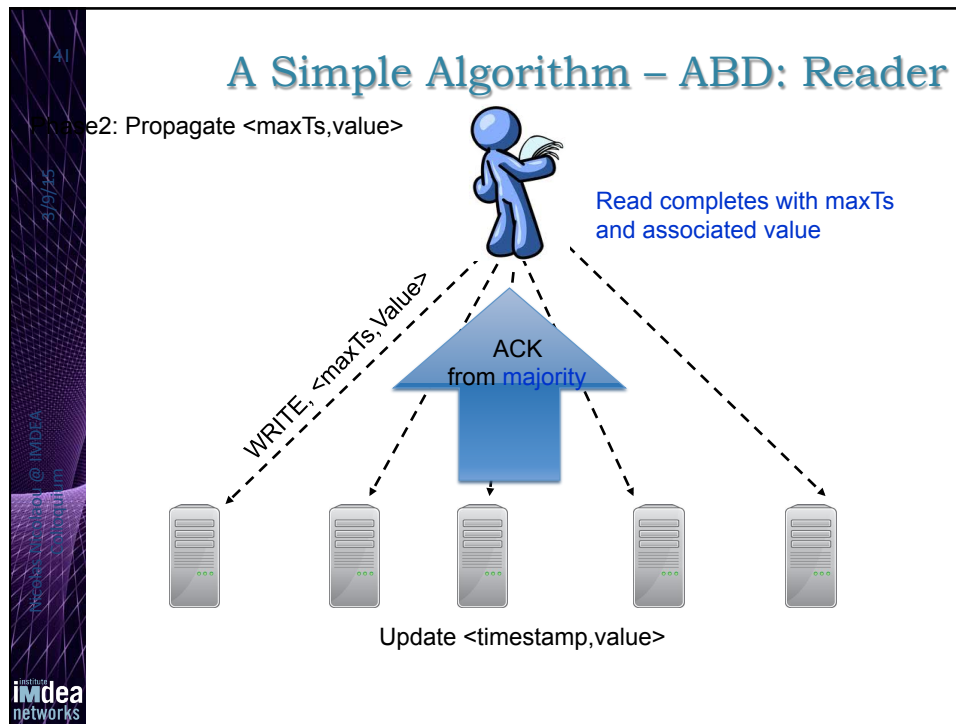
40

A Simple Algorithm – ABD: Reader

Case 1: Discover maximum timestamp



imdea
networks



42

Tagging the values

- TAG : $\langle \text{timestamp}, \text{wid} \rangle$ pair
 - tag1 > tag2 if either:
 - tag1.timestamp > tag2.timestamp, or
 - tag1.timestamp = tag2.timestamp AND tag1.wid > tag2.wid
- Why wid is necessary?
 - Separate writes with the same timestamp

imdea networks