



Trade-offs in Implementing Atomic  
Multi-Writer, Multi-Reader Registers in  
Asynchronous Message-passing Systems

Chryssis Georgiou  
**University of Cyprus**

May 17, 2011

This work is partially funded by the Cyprus Research Promotion Foundation and  
co-funded by the Republic of Cyprus and the European Regional Development Fund

17/5/2011 – Chryssis Georgiou ©

## Motivation for this Research

17/5/2011 – Chryssis Georgiou ©

## Survivability of Data

- ▶ Data Survivability is **essential** in today's systems and applications
- ▶ Popular approach – RAID [Patterson, Gibson, Katz 1988]
  - ▶ Use of a Redundant Array of Inexpensive Disks
  - ▶ Mirroring or erasure code is used to prevent loss of data upon a disk failure
  - ▶ A RAID system contains
    - ▶ a single box
    - ▶ residing at a single physical location
    - ▶ connected to clients via a single network interface

**Single point of failure**

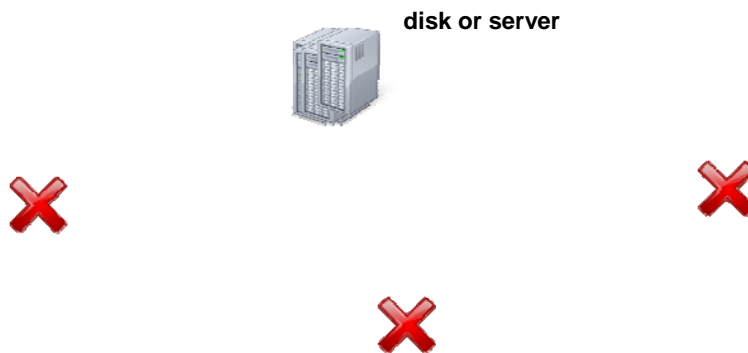
[Chockler, Keidar, Guerraoui, Vukolic 09]  
[Schroeder, Gibson 07]

▶ 3

17/5/2011 – Chryssis Georgiou ©

## Distributed Storage

- ▶ Uses Replication



- ▶ Fault-Tolerance, Availability and Geographical Proximity

▶ 4

17/5/2011 – Chryssis Georgiou ©

## No Free Lunch

---

### **Great Challenge:**

Maintain **consistency** among the replicas  
despite system **asynchrony** and **failures**  
And do so **efficiently**

#### **Akamai's Top Ten Challenges**

Harald Prokop, Senior Vice President of Engineering, July 2008

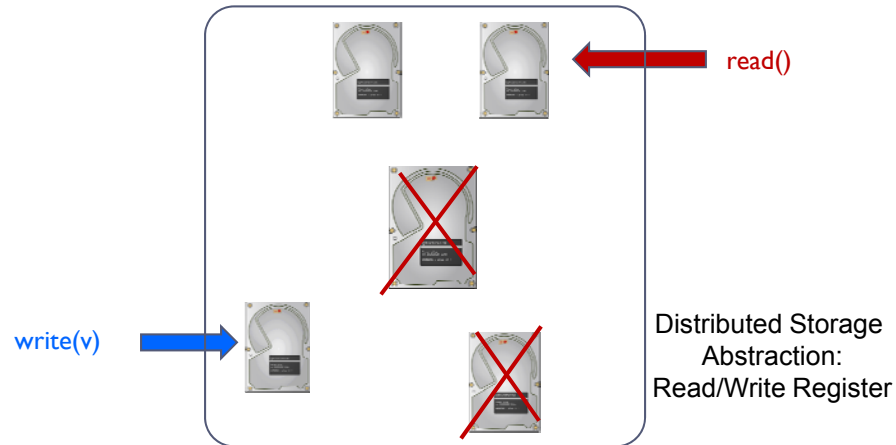
▶ 5

17/5/2011 – Chrysis Georgiou ©

## Problem Abstraction

17/5/2011 – Chrysis Georgiou ©

## Read/Write Object (register)



- ▶ Read/Write operations – invocation and response
- ▶ Concurrency
- ▶ Consistency Semantics

▶ 7

17/5/2011 – Chryssis Georgiou ©

## Operation Relations

- ▶ Precedence Relations for two operations  $\pi_1, \pi_2$ :

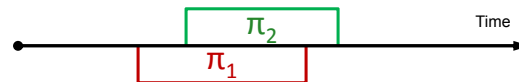
- ▶  $\pi_1$  **precedes**  $\pi_2$  if the response of  $\pi_1$  happens **before** the invocation of  $\pi_2$



- ▶  $\pi_1$  **succeeds**  $\pi_2$  if the invocation of  $\pi_1$  happens **after** the response of  $\pi_2$



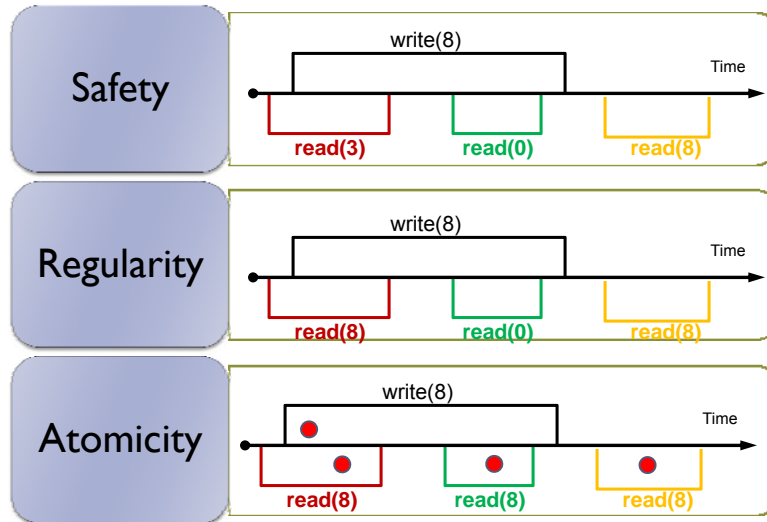
- ▶  $\pi_1$  is **concurrent** with  $\pi_2$  if  $\pi_1$  **neither precedes nor succeeds**  $\pi_2$



▶ 8

17/5/2011 – Chryssis Georgiou ©

## Consistency Semantics [Lamport 1986]



▶ 9

17/5/2011 – Chrysis Georgiou ©

## Atomic Read/Write Register [Lamport 1986]

- ▶ An **atomic Read/Write register** allows concurrent processes
  - ▶ to **share information** through a common variable
  - ▶ **as if they were** accessing this variable in a sequential manner
- ▶ It supports **read** and **write** operations
  - ▶ Read() returns the latest written value of the object
  - ▶ Write(v) sets the value to v and returns OK

▶ 10

17/5/2011 – Chrysis Georgiou ©

## Applicability

- ▶ The atomic register is **fundamental** in distributed computing and is at the heart of a large number of distributed algorithms [Attiya & Welch 1998/2004]
- ▶ Atomic register implementations are used as building blocks for complex distributed storage systems [Fleet 2000, SBQ-L 2002, PASIS 2004, HP's FAB 2006, Amazon's Dynamo 2007]
- ▶ The can also be used directly to build distributed file systems [Chockler, Keidar, Guerraoui, Vukolic 09]

▶ 11

17/5/2011 – Chrysis Georgiou ©

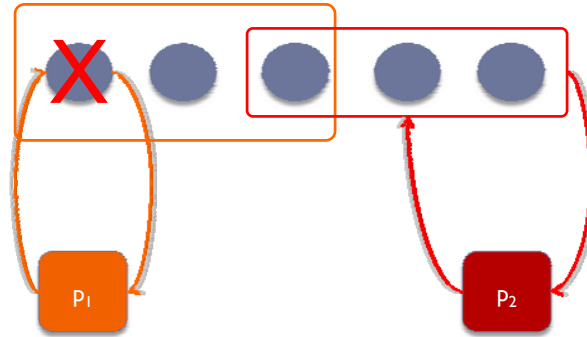
## Challenge: Ordering of read/write operations

- ▶ Using the “time” at which each operation is invoked
  - ▶ **Clock Synchronization**
  - ▶ **Access to centralize clock** (e.g. , GPS)
- ▶ Associate a sequence number with each value written
  - ▶  $tags = \langle timestamp, wid \rangle$ , [Attiya, Bar-Noy, Dolev 1995]
    - ▶ Timestamp: counter (positive integer)
    - ▶ wid: writer id
  - ▶  $tag1 > tag2$  if either:
    - ▶  $tag1.timestamp > tag2.timestamp$ , or
    - ▶  $tag1.timestamp = tag2.timestamp$  AND  $tag1.wid > tag2.wid$

▶ 12

17/5/2011 – Chrysis Georgiou ©

## Challenge: Communication Rounds

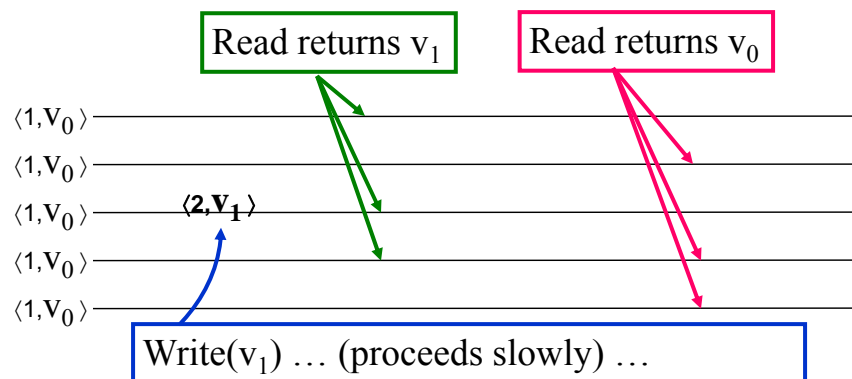


▶ 13

17/5/2011 – Chrysis Georgiou ©

## Multiple Round-Trips

- ▶ Consider the following example [Attiya, Bar-Noy, Dolev 1995]:



▶ 14

17/5/2011 – Chrysis Georgiou ©

## Complexity Measure-Operation Latency

Operation Latency is measured in  
**Communication Rounds (round-trips)**

A process  $p$  performs a **communication round** for an operation  $\pi$  if

- Sends a message  $m$  regarding  $\pi$  to a subset (potentially all) of processes
- Any process that receives  $m$ , **replies** to  $p$
- Process  $p$  **collects** “enough” of such replies and proceeds accordingly.

▶ 15

17/5/2011 – Chryssis Georgiou ©

## Research Questions

What is the **operation latency** of atomic register implementations in a **fail-prone, message-passing, asynchronous** distributed system?

What are the **trade-offs** to achieve such performance?

▶ 16

17/5/2011 – Chryssis Georgiou ©



## Background and Prior Work

17/5/2011 – Chrysis Georgiou ©

### Model of Computation

- ▶ **Asynchronous, Message-Passing model**
  - ▶ Process sets: writers  $W$ , readers  $R$ , servers  $S$  (replica hosts)
  - ▶ Reliable Communication Channels
  - ▶ Well Formedness
- ▶ **Environments:**
  - ▶ Single Writer, Multiple Readers:  $|W|=1, |R|\geq 1$
  - ▶ Multiple Writers, Multiple Readers:  $|W|\geq 1, |R|\geq 1$
- ▶ **Failures:**
  - ▶ Crash Failures
    - ▶ A subset of the servers may crash, all writer and reader processes may crash
- ▶ **Correctness: Atomicity (safety), Termination (liveness)**

▶ 18

17/5/2011 – Chrysis Georgiou ©

## Fastness

- ▶ **Fast Operation**
  - ▶ Completes at the end of its first round
- ▶ **Fast Implementation of an atomic register**
  - ▶ All operations are fast
- ▶ **Communication scheme**
  - ▶ Message delivery: Servers ↔ Clients
  - ▶ No server to server or client to client communication

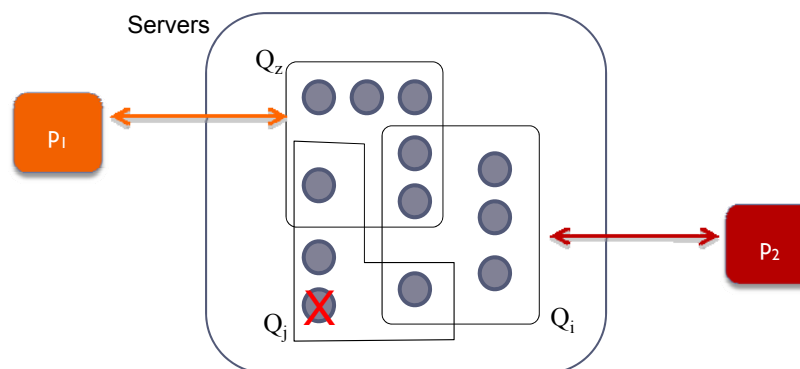
Data-centric communication

[Martin, Alvisi, Dahlin 2002]  
[Chockler, Keidar, Guerraoui, Vukolic 2009]

▶ 19

17/5/2011 – Chrysis Georgiou ©

## Tools: Quorum systems

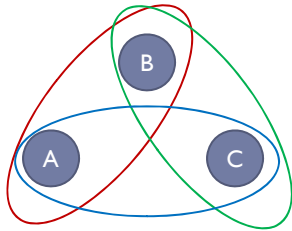


- ▶  $Q_i, Q_j, Q_z$  are **quorums**
- ▶ **Quorum System** is the set  $\{Q_i, Q_j, Q_z\}$ 
  - ▶ Property: every pair of quorums intersects
- ▶ Every R/W operation communicates with a single quorum
- ▶ **Faulty Quorum**: Contains a faulty process

▶ 20

17/5/2011 – Chrysis Georgiou ©

## Quorum Systems Examples



Majorities [Thomas79,Gifford79]

	■		■
■	■	■	■
	■		■
■	■	■	■

Matrix Quorums [Vitanyi, Awerbuch 1992]

▶ 21

17/5/2011 – Chrysis Georgiou ©

## Prior Work: Traditional Implementations



- e.g., [Attiya, Bar-Noy, Dolev 95]
- Single round writes
- Two round reads
  - Phase 1: Obtain latest value (max tag)
  - Phase 2: Propagate latest value (max tag)
    - Folklore belief: “Reads must Write”



- e.g., [Lynch Shvartsman 97, 02, Englert, Shvartsman 00]
- Two round writes
  - Phase 1: Discover latest value (max tag)
  - Phase 2: Order new value after the latest and propagate
    - Belief: “Writes must Read”
- Two round reads

▶ 22

17/5/2011 – Chrysis Georgiou ©

## Algorithm: Simple

### Write Protocol: two rounds

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

### Read Protocol: two rounds

- P1: Query a single quorum for the latest tag
- P2: Propagate  $\langle \text{maxtag}, v \rangle$  to a single quorum

### Server Protocol: passive role

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

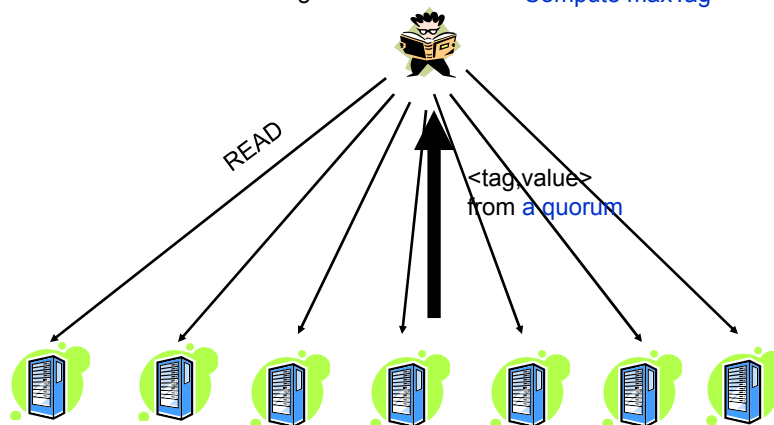
▶ 23

17/5/2011 – Chryssis Georgiou ©

## Reader

Round1: Discover maximum tag

Compute maxTag

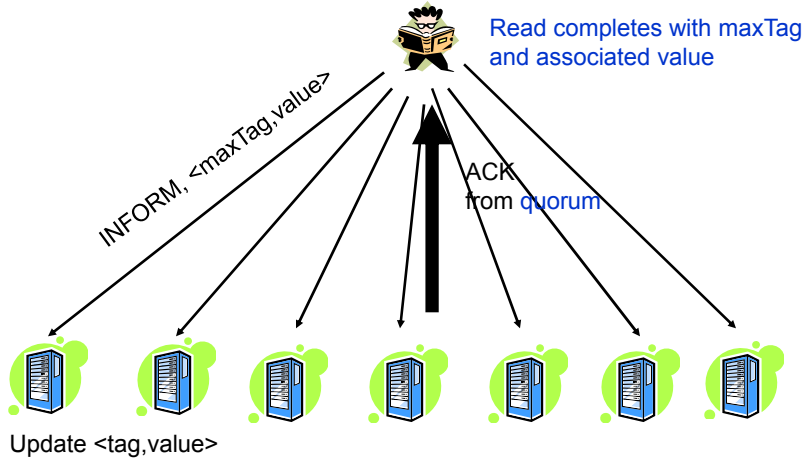


▶ 24

17/5/2011 – Chryssis Georgiou ©

# Reader

Round2: Propagate  $\langle \text{maxTag}, \text{value} \rangle$

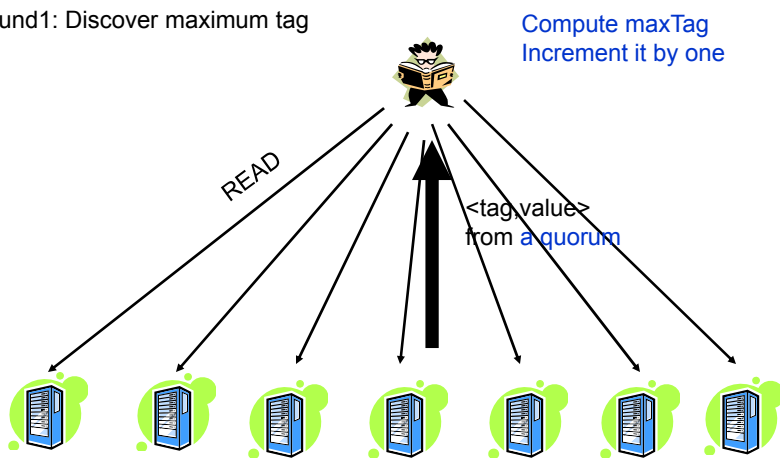


▶ 25

17/5/2011 – Chryssis Georgiou ©

# Writer

Round1: Discover maximum tag

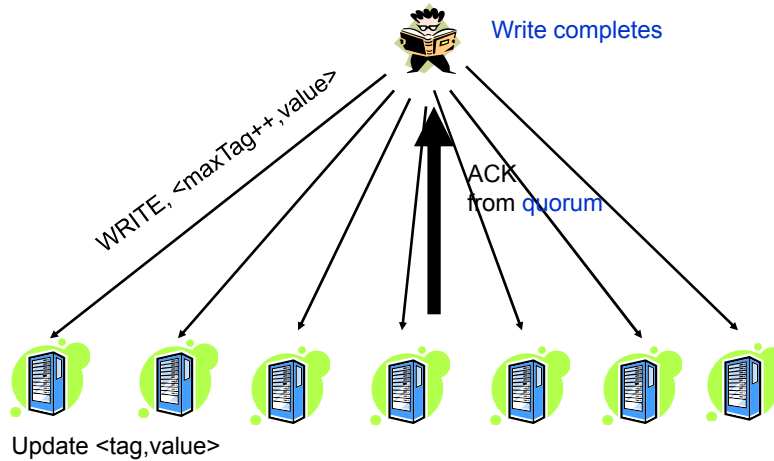


▶ 26

17/5/2011 – Chryssis Georgiou ©

## Writer

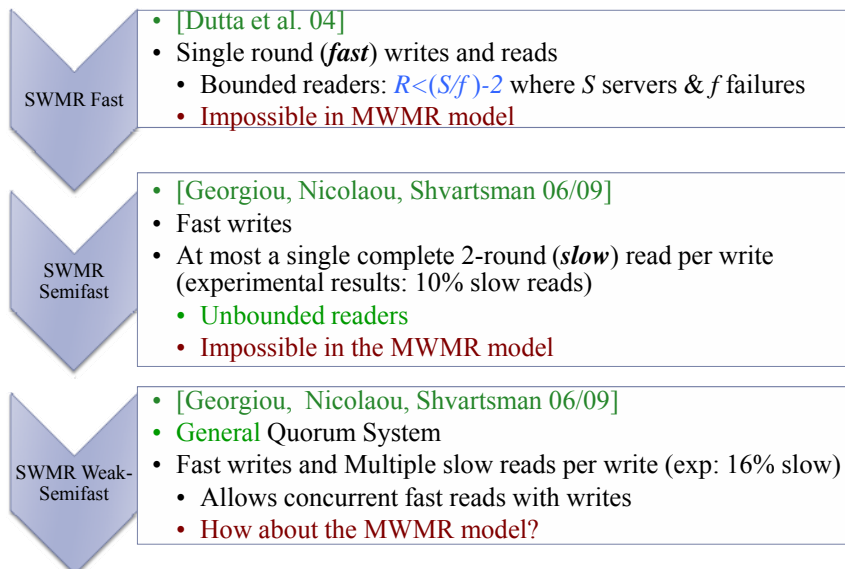
Round2: Propagate  $\langle \text{maxTag}, \text{value} \rangle$



▶ 27

17/5/2011 – Chryssis Georgiou ©

## Prior Work: Fast Operations



▶ 28

17/5/2011 – Chryssis Georgiou ©

## Fastness in the MWMR Setting

17/5/2011 – Chryssis Georgiou ©

### Question

---

Can we allow **some** reads to be fast (single round reads) and still guarantee atomicity in the MWMR settings?

Answer: **YES!!**

▶ 30

17/5/2011 – Chryssis Georgiou ©

## Tool: Quorum Views

### Idea:

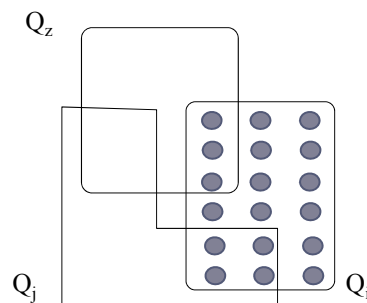
- ▶ Try to determine the state of the write operation based on the distribution of the latest value in the replied quorum.
- ▶ Write State in the First Round of Read Operation
  - Determinable => Read is **Fast**
  - Undeterminable => Read is **Slow**

▶ 31

17/5/2011 – Chryssis Georgiou ©

## Determinable Write - Qview(1)

- ▶ All members of a quorum contain the maxTag



(Potentially) Write Completed

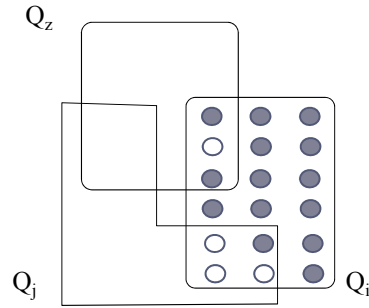
▶ 32

17/5/2011 – Chryssis Georgiou ©



## Determinable Write - Qview(2)

- ▶ Every intersection contains a member with tag < maxTag



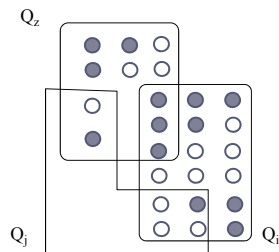
(Definitely) Write <maxTag,v> Incomplete

▶ 33

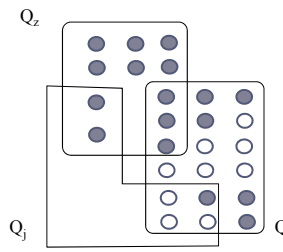
17/5/2011 – Chrysis Georgiou ©

## Undeterminable Write - Qview(3)

- ▶ There is intersection with all its members with tag = maxTag



qV(3) and Incomplete Write



qV(3) and Complete Write

Undeterminable => second Com. Round

▶ 34

17/5/2011 – Chrysis Georgiou ©

## 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<sub>0</sub>

### Server Protocol: passive role

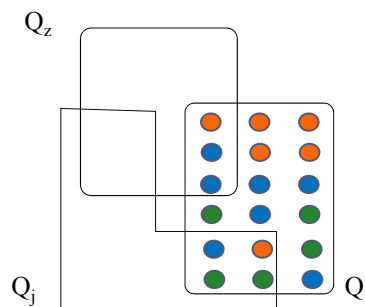
- Receive requests, update local timestamp and return  $\langle \text{tag}, v \rangle$

▶ 35

17/5/2011 – Chrysis Georgiou ©

## Quorum Views and Multiple Writers

- ▶ 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  $\text{tag1} < \text{tag2} < \text{tag3}$

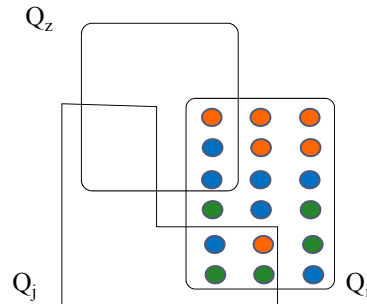


▶ 36

17/5/2011 – Chrysis Georgiou ©

## 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**

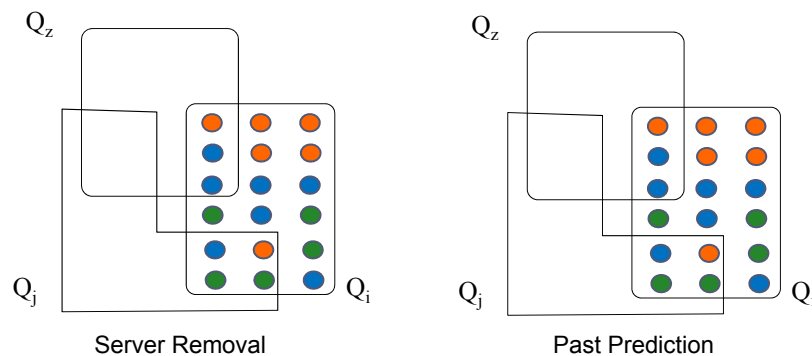


▶ 37

17/5/2011 – Chrysis Georgiou ©

## 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  $Q_i$** 
    - ▶ **Qview(1)**: all remaining servers contain  $\langle \text{tag1}, v1 \rangle$

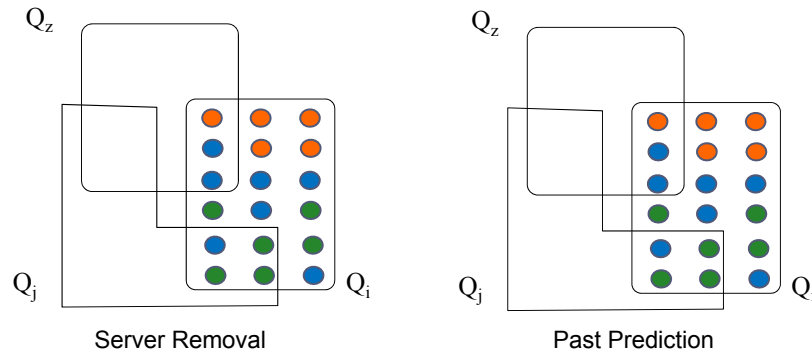


▶ 38

17/5/2011 – Chrysis Georgiou ©

## 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)



▶ 39

17/5/2011 – Chryssis Georgiou ©

## What about writes?

Can we devise **MWMMR atomic register implementations** that allow executions that contain both **fast read and write operations**?

**Answer: YES!!**

▶ 40

17/5/2011 – Chryssis Georgiou ©

## New Technique – SSO [Englert et. al OPODIS 2009]

### ▶ 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

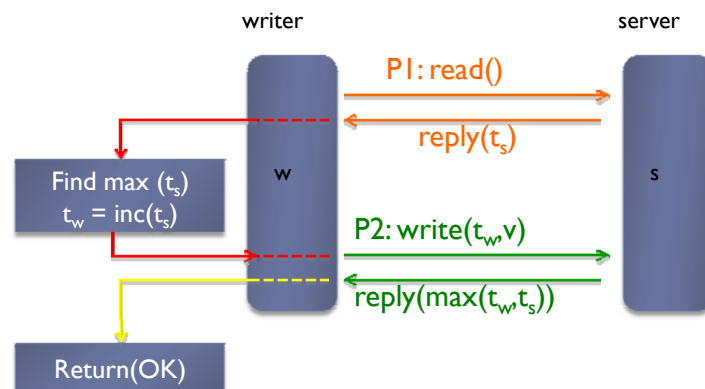
### ▶ SFW Algorithm

- ▶ Deploys the SSO technique
- ▶ Enables **Fast Writes and Reads** -- first such algorithm

▶ 41

17/5/2011 – Chrysis Georgiou ©

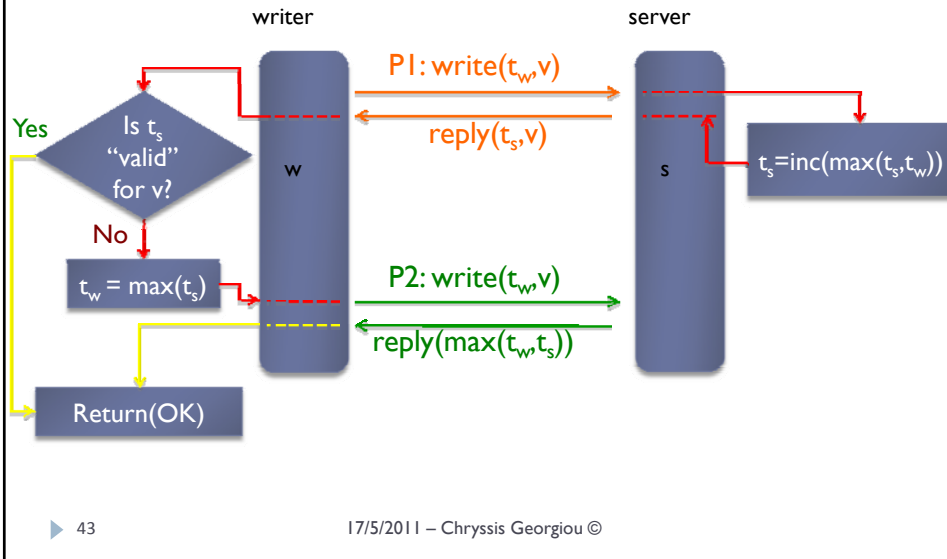
## Traditional Writer-Server Interaction



▶ 42

17/5/2011 – Chrysis Georgiou ©

## SSO Writer-Server Interaction



## Definition: $n$ -wise Quorum Systems

**Definition:** A quorum system  $Q$  is an  $n$ -wise quorum system, if:

$$Q = \{Q : Q \subseteq S\} \text{ where } \forall A \subseteq Q : |A| = n \text{ and } \bigcap_{Q \in A} Q \neq \emptyset$$

## Algorithm: SFW

### Write Protocol: one or two rounds

- P1: send  $v$  and gather **candidate** tags from a quorum
- Exists tag  $t$  propagated in a **bigger** than  $(n/2-1)$ -wise intersection (PREDICATE)
  - **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)
  - **YES** – return the value written by that write => **FAST**
  - **NO** - propagate the largest confirmed tag to a quorum => **SLOW**

### Server Protocol

- **Increment tag** when receive write request and **record the latest writes**
- Upon read/write request **send the recording set**

▶ 45

17/5/2011 – Chrysis Georgiou ©

## Predicates: Read and Write

**Writer predicate for a write  $\omega$  (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  $\rho$  (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.  $(\bigcap_{Q \in Q^j \cup \{Q\}} Q) \subseteq MS$ .

▶ 46

17/5/2011 – Chrysis Georgiou ©

## Lower bound

---

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

▶ 47

17/5/2011 – Chrysis Georgiou ©

## Remark

---

**Remark:** Algorithm SFW is near optimal since it allows up to  $\left(\frac{N}{2}-1\right)$  consecutive, quorum shifting, fast writes.

▶ 48

17/5/2011 – Chrysis Georgiou ©



## Trade-offs in MWMR

Trades Write  
Speed for Fast  
Reads

- Quorum Views in the MWMR environment
- Algorithm CWFR
- Traditional two round writes
- Some single round reads – even when reads are concurrent with writes
- Utilizes any General Quorum System

Trades Quorum  
Generality for Fast  
Writes

- Algorithm SFW
- Server Side Ordering
- Allows both single round reads and writes in MWMR
- Fastness depends on  $n$ -wise quorum intersections
- $n-1$  consecutive fast writes are possible in MWMR
- SFW near optimal – Allows  $n/2$  consecutive fast writes

▶ 49

17/5/2011 – Chrysis Georgiou ©

## Recent Advancements

17/5/2011 – Chrysis Georgiou ©

## Problem: Communication vs Computation

**Theorem:** The computation of the read and write predicates used by algorithm SWF is **NP-complete**.

We provide a reduction from 3-SAT

▶ 51

17/5/2011 – Chrysis Georgiou ©

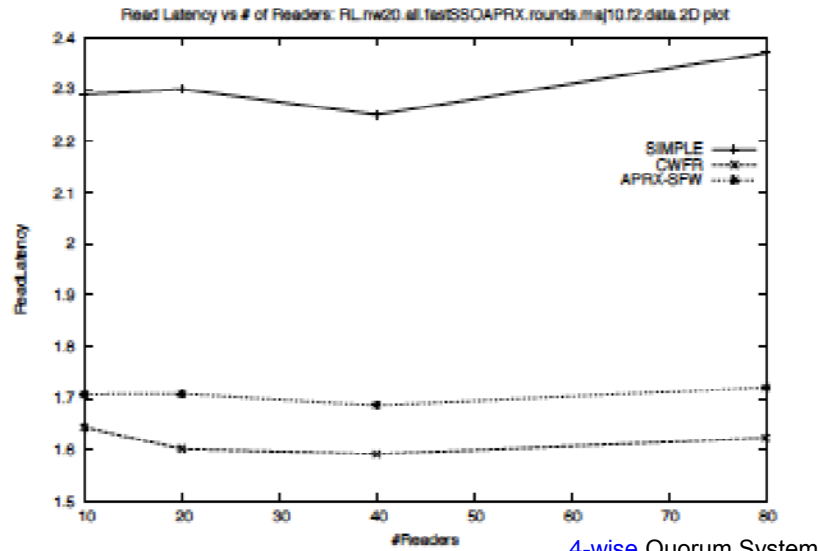
## Approximation Algorithm: AprxSFW

- ▶ Greedy algorithm
  - ▶ Uses Set Cover greedy approximation algorithm at its core for the predicate computation
- ▶ Log-Approximation
  - ▶ Invalidates RP and WP in a factor of  $\log|S|$  times
  - ▶ A factor of  $\log|S|$  more second communication rounds
  - ▶ Does not affect correctness
- ▶ Preliminary empirical evaluation using the NS2 simulator
  - ▶ Ongoing work: Evaluation on Planetlab

▶ 52

17/5/2011 – Chrysis Georgiou ©

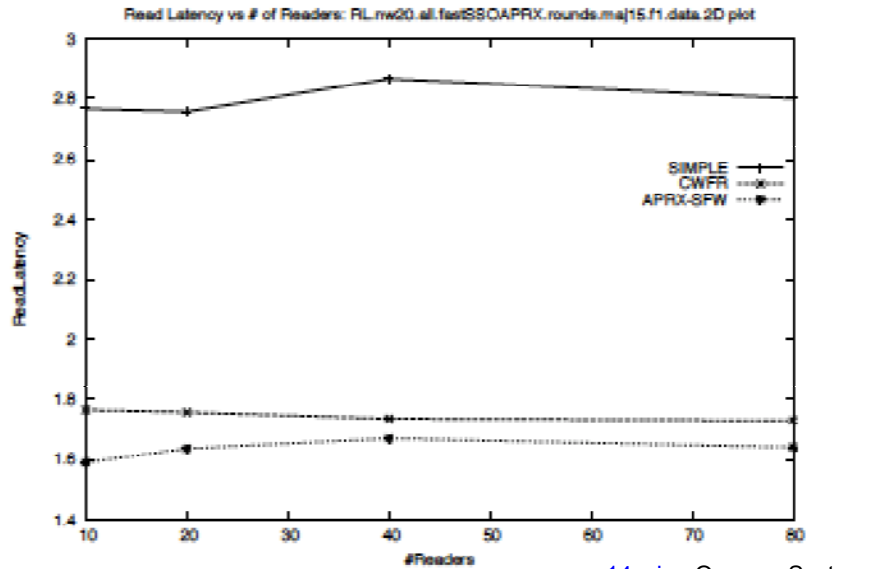
### AprxSFW – CWFR Empirical Evaluation



► 53

17/5/2011 – Chrysis Georgiou ©

### AprxSFW – CWFR Empirical Evaluation



► 54

17/5/2011 – Chrysis Georgiou ©



**Gracias!**