

Quorum Views: Client Side Tools for Atomic Distributed Register Implementations

Nicolas Nicolaou
 University of Connecticut

Funded by the Cyprus Research Promotion Foundation, the Republic of Cyprus and the European Regional Development Fund

What is a Distributed Storage System?

Distributed Storage Abstraction

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

▶ 2 Nicolas Nicolaou 17/3/2011

Consistency Semantics [Lamport86]

Safety: write(8), read(3), read(0), read(8)
 Regularity: write(8), read(8), read(0), read(8)
 Atomicity: write(8), read(8), read(8), read(8)

▶ 3 Nicolas Nicolaou 17/3/2011

Complexity Measure-Operation Latency

- ▶ Consistent Register Implementations
 - ▶ **Message-Passing, Asynchronous** model
 - ▶ Access **multiple replicas** per operation
 - ▶ Perform **multiple accesses** per operation

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

▶ 4 Nicolas Nicolaou 17/3/2011

Fastness of Atomic Implementations

Traditional SWMR

- Single round writes
- Two round reads
- Phase 1: Obtain latest value
- Phase 2: Propagate latest value
- Folklore belief: "Reads must Write"

SWMR Fast

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

SWMR Semifast

- **Fast 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**

▶ 5 Nicolas Nicolaou 17/3/2011

Quorum Systems

- ▶ **Quorum System:**
A collection of sets with non-empty pairwise intersections
- ▶ Numerous results exploit quorums for
 - ▶ MWMR and Dynamic Atomic Memory Implementations
 - ▶ "Slow" operations
 - ▶ Some achieve fast reads when a write is confirmed [DGLSW03, CGGMS05]
- ▶ Refined Quorum Systems - Guerraoui and Vukolic [PODC2007]
 - ▶ Introduce Fast Operations in Quorum-Based systems
 - ▶ **Rely on operation timeout**

▶ 6 Nicolas Nicolaou 17/3/2011

Question

Can we obtain Fast/Semifast Implementations in a general quorum-based framework?

▶ 7

Nicolas Nicolaou

17/3/2011

First the Bad News

- ▶ Quorum-Based Fast Implementations are **not Fault-Tolerant**
 - ▶ Possible iff there exists a **common intersection** among all the quorums
 - ▶ A **single failure** in the common intersection collapses the Quorum system.
- ▶ Quorum-Based Semifast Implementations are **not Fault-Tolerant**
 - ▶ One of the properties of the semifast definition is violated if no common intersection exists among all the quorums.
 - ▶ A single complete "slow" read operation is not enough for every write operation

▶ 8

Nicolas Nicolaou

17/3/2011

Now the Good News

- ▶ Trade speed for efficiency and fault-tolerance
 - ▶ Allow multiple "slow" reads per write operation but maintain the fast behavior when possible
- ▶ To do so, we introduce **Quorum Views**

▶ 9

Nicolas Nicolaou

17/3/2011

Model: Definitions

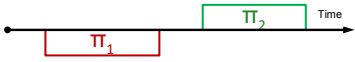

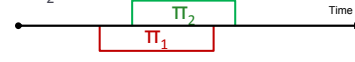
- ▶ Asynchronous, Message-Passing model
 - ▶ Process sets: writers W, readers R, servers S (replica hosts)
 - ▶ Reliable Communication Channels
 - ▶ Well Formedness
- ▶ Environments:
 - ▶ SWMR: $|W|=1, |R|\geq 1$
 - ▶ MWMR: $|W|\geq 1, |R|\geq 1$
- ▶ Failures:
 - ▶ Crash Failures
- ▶ Correctness: Atomicity (safety), Termination (liveness)
- ▶ Order write operations
 - ▶ SWMR: $\langle \text{timestamp}, \text{value} \rangle$ pairs
 - ▶ MWMR: $\langle \text{tag}, \text{value} \rangle$ pairs

▶ 10

Nicolas Nicolaou

17/3/2011

Definition: Operation Relations

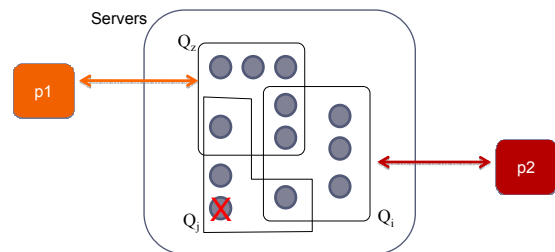
- ▶ Precedence Relations for two operations π_1, π_2 :
 - ▶ π_1 **precedes** π_2 if the response of π_1 happens **before** the invocation of π_2

 - ▶ 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


▶ 11

Nicolas Nicolaou

17/3/2011

Definition: Quorum systems



- ▶ Q_1, Q_2, Q_3 are **quorums**
- ▶ **Quorum System** is the set $\{Q_1, Q_2, Q_3\}$
 - ▶ Property: every pair of quorums intersects
- ▶ Every R/W operation communicates with a single quorum
- ▶ **Faulty Quorum**: Contains a faulty process

▶ 12

Nicolas Nicolaou

17/3/2011

Definition: Fastness

- ▶ A process p performs a **communication round** during an operation π if:
 - ▶ p sends a message m to a set of servers for π
 - ▶ Any server that receives m replies to p
 - ▶ Once p receives responses from a single quorum completes π or proceeds to a next communication round
- ▶ **Fast Operation**
 - ▶ Completes at the end of its first round
- ▶ **Fast Implementation**
 - ▶ All operations are fast
- ▶ **Communication scheme**
 - ▶ Message delivery: **Servers to Clients**
 - ▶ No server to server or client to client communication

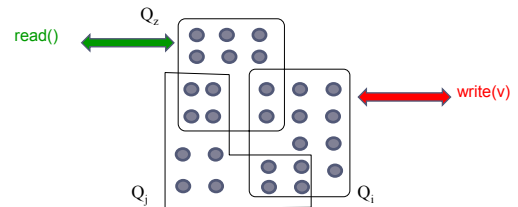
▶ 13

Nicolas Nicolaou

17/3/2011

Non-Robust Fast Implementations: Proof Sketch

- ▶ **Execution a:**
 - ▶ Complete **write(v)** $\Rightarrow Q_i$
 - ▶ Complete **read()** $\Rightarrow Q_z$
 - ▶ **read()** returns v to preserve atomicity



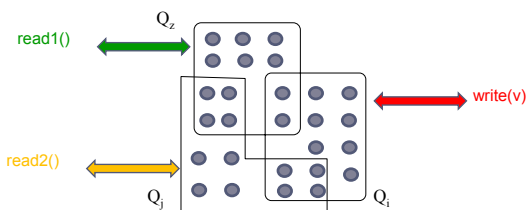
▶ 14

Nicolas Nicolaou

17/3/2011

Not Robust Fast Implementations (Proof Sketch)

- ▶ **Execution b:**
 - ▶ Incomplete **write(v)** $\Rightarrow Q_i \cap Q_z$
 - ▶ Complete **read1()** $\Rightarrow Q_z$
 - ▶ **read1()** cannot distinguish between executions a and b, thus returns v
 - ▶ Complete **read2()** $\Rightarrow Q_i$
 - ▶ **read2()** returns an **older** value since does not observe v



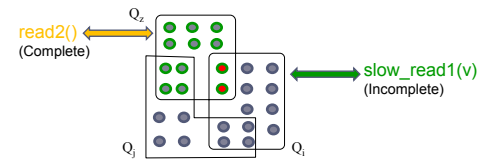
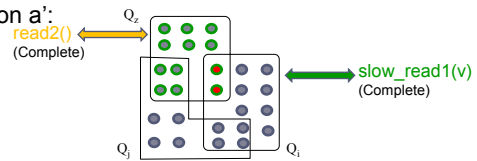
▶ 15

Nicolas Nicolaou

17/3/2011

Not Robust SemiFast Implementations (Proof Sketch)

- ▶ **Execution a':**
 - ▶ Complete **read2()** (Complete) $\Rightarrow Q_z$
 - ▶ Complete **slow_read1(v)** (Complete) $\Rightarrow Q_i$
- ▶ **Execution b':**
 - ▶ Complete **read2()** (Complete) $\Rightarrow Q_z$
 - ▶ Complete **slow_read1(v)** (Incomplete) $\Rightarrow Q_i$



▶ 16

Nicolas Nicolaou

17/3/2011

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 \Rightarrow Read is **Fast**

Undeterminable \Rightarrow Read is **Slow**

▶ 17

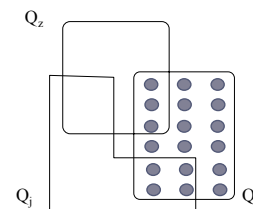
Nicolas Nicolaou

17/3/2011

Determinable Write - Qview(1)

- ▶ All members of a quorum contain $\max T_s$

$$[qView(1)]: \forall s \in Q_i : s.ts = \max TS$$



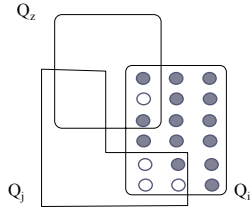
▶ 18

Nicolas Nicolaou

17/3/2011

Determinable Write - Qview(2)

- ▶ Every intersection contains a member with $ts < \max TS$
 $[qView(2)]: \forall j \neq i, \exists A \subseteq Q_i \cap Q_j, s.t. A \neq \emptyset \text{ and } \forall s \in A: s.ts < \max TS$



(Definitely) Write $\langle \max Tag, v \rangle$ Incomplete

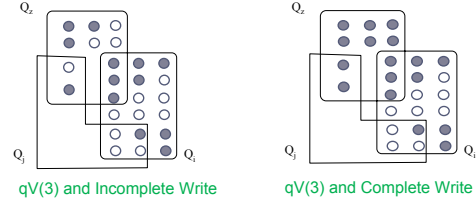
▶ 19

Nicolas Nicolaou

17/3/2011

Undeterminable Write - Qview(3)

- ▶ There is intersection with all its members with $ts = \max TS$
 $[qView(3)]: \exists s \in Q_i, s.ts = \max TS \text{ and } \exists j \neq i, \forall s \in Q_i \cap Q_j: s.ts = \max TS$



Undeterminable => second Com. Round

▶ 20

Nicolas Nicolaou

17/3/2011

Algorithm: SLIQ

Write Protocol: one round

- P1: Writer increments ts and propagates the $\langle ts, v \rangle$ to a quorum

Read Protocol: one or two rounds

- P1: send read requests and wait for replies from a quorum Q
 - $QView_Q(1)$ – **Fast** and return $\max TS$
 - $QView_Q(2)$ – **Fast** and return $\max TS - 1$
 - $QView_Q(3)$ – **Slow** proceed to P2 and return $\max TS$
- P2: propagate $\langle \max TS, v \rangle$ to a quorum and return $\langle \max TS, v \rangle$

Server Protocol: passive role

- Receive requests, update local timestamp and return $\langle ts, v \rangle$

▶ 21

Nicolas Nicolaou

17/3/2011

SLIQ: Simulation Results

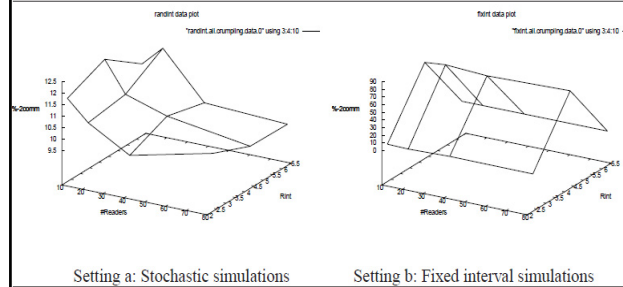


Fig. 3. Simple runs using Crumbling Walls

▶ 22

Nicolas Nicolaou

17/3/2011

Follow-Up Question

Can Quorum Views be used to obtain fast operations in the **MWMM** environment?

▶ 23

Nicolas Nicolaou

17/3/2011

Tagging the values

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

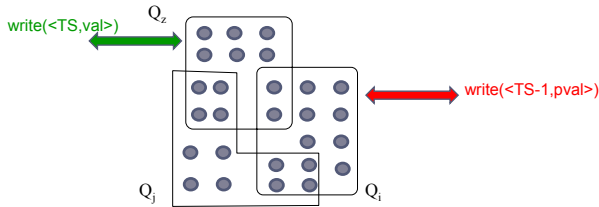
▶ 24

Nicolas Nicolaou

17/3/2011

Why be fast in Qview(2) in SWMR?

- Single writer
 - If it writes $\langle TS, val \rangle \Rightarrow$ it completed writing $\langle TS-1, pval \rangle$



Any read operation will observe either QView1 or QView3 for $\langle TS-1, pval \rangle$

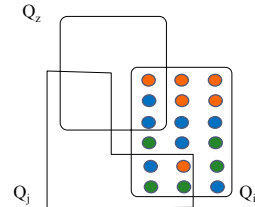
▶ 25

Nicolas Nicolaou

17/3/2011

What happens in MWMR?

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



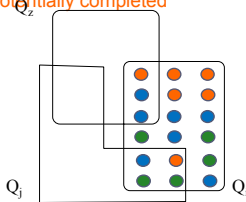
▶ 26

Nicolas Nicolaou

17/3/2011

Idea: Uncover the Past

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



▶ 27

Nicolas Nicolaou

17/3/2011

Algorithm: CWFR

Traditional Write Protocol: two rounds

- P1: Query a single quorum for the latest tag
- P2: Increment the max tag, send $\langle newtag, v \rangle$ 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 maxTag₀

Server Protocol: passive role

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

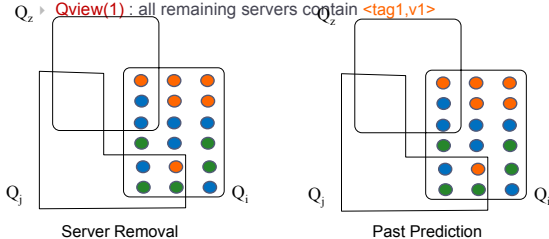
▶ 28

Nicolas Nicolaou

17/3/2011

Read Iteration: Discard Incomplete Tags

- For values $\langle tag1, v1 \rangle$, $\langle tag2, v2 \rangle$, $\langle tag3, v3 \rangle$:
 - $\langle tag3, v3 \rangle$ not completed: remove servers that contain $\langle tag3, v3 \rangle$
 - $\langle tag2, v2 \rangle$ not completed: remove servers that contain $\langle tag2, v2 \rangle$
 - $\langle tag1, v1 \rangle$ potentially completed in Q_i
- $Q_z \rightarrow QView(1)$: all remaining servers contain $\langle tag1, v1 \rangle$



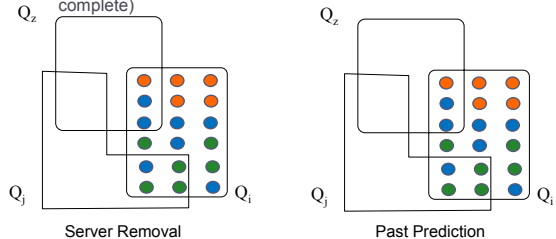
▶ 29

Nicolas Nicolaou

17/3/2011

Read Iteration: Discard Incomplete Tags

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



▶ 30

Nicolas Nicolaou

17/3/2011

CWFR: Empirical Results

- ▶ In progress experimentation on Planetlab
 - ▶ Planetary scale real-time environments
- ▶ Measure the percentage of "slow" reads
- ▶ Get the average of operation latency
- ▶ Compare with classic 2 round approaches

Closing Remarks

