

# Seeking Fastness in Multiple-Writer Multiple-Reader Atomic Register Implementations

Π Ε Ν Ε Κ / 0609/31



Chryssis Georgiou, Nicolas Nicolaou, Alexander Russell, Alexander Shvartsman

## Problem - Background

**What is a Distributed Storage System?**

**Traditional Implementations**

- Traditional SWMR**
  - Single round writes
  - Two round reads
  - Phase 1: Obtain latest value
  - Phase 2: Propagate latest value
  - Folklore belief: "Reads must Write"
- Traditional MWMR**
  - Two round writes
  - Phase 1: Discover latest value
  - Phase 2: Order new value after the latest and propagate
  - Belief: "Writes must Read"
  - Two round reads

**The Era of Fast Implementations...**

- Fast SWMR**
  - Single round (Fast) Writes
  - All fast reads with bounded readers [Dutta et al. 04]
  - Semi-fast: A single slow read per write and unbounded readers [Georgiou et al. 06]
  - Not applicable to the MWMR problem
- Fast MWMR**
  - Algorithm SFW [Englert et al. 2009]
  - Server Side Ordering (SSO)
  - Order values at the servers
  - Avoid discovery of the latest value from the writer
  - Enables Fast (Single Round) Writes and Reads -- first such algorithm

## Distributed Setting

**Model**

**Environment:**

- Asynchronous
- Message-Passing
- MWMR:  $|W| \geq 1, |R| \geq 1$

**Process sets:**

- writers  $W$ ,
- readers  $R$ ,
- servers  $S$  (replica hosts)

**Failures: Crash Failures**

**Correctness:**

- Atomicity (safety)
- Termination (liveness)

**Definition: 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  $\pi_2$

**Consistency Semantics [Lamport86]**

- Safety**
- Regularity**
- Atomicity**

**Complexity Measure**

Communication Delays (round-trips) + Computation = Operation Latency

**Definition: Quorum systems**

Servers  $S$  are partitioned into quorums  $Q_1, Q_2, \dots, Q_k$ .

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

## Motivation

**The Weak Side of SFW**

- Predicates are Computationally Hard
  - NP-Complete (Shown in this paper)
- Restriction on the Quorum System
  - Deploys  $n$ -wise Quorum Systems (every  $n$  quorums intersect)
  - Guarantees fastness iff  $n > 3$

**The Good News...**

- Approximation Algorithm (APRX-SFW)
  - Polynomial
  - Log-approximation
    - $\log|S|$  times the optimal number of fast operations
- Algorithm CWFR
  - Based on Quorum Views
  - SWMR prediction tools
  - Fast operations in General Quorum Systems
  - Trades Speed of Write operations
  - Two Round Writes

## Algorithm SFW

**Key Idea**

**SSO: Server Side Ordering**

- Tag is incremented by the server and not by the writer.
- Generated tags may be different across servers
- Clients decide operation ordering based on server responses

**SSO Algorithm**

- Enables Fast Writes and Reads -- first such algorithm
- Allows Unbounded Participation

**Algorithm: SFW**

**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  $(n/2-1)$ -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

**SSO Writer-Server Interaction**

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

**Theorem: It is impossible to get MWMR safe register implementations that exploit an  $N$ -wise quorum system, if  $|W \cup R| > N-1$**

**Remark: SSO algorithm is near optimal since it allows up to  $\lfloor \frac{N}{2} \rfloor$  consecutive, quorum shifting, fast writes.**

## Algorithm SIMPLE

**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  $\langle \text{server.tag}, v \rangle$

## Algorithm APRX-SFW

**K-SET-INTERSECTION: (captures both PR and PW)**

Given a set of elements  $U$ , a subset of those elements  $M \subseteq U$ , a set of subsets  $Q = \{Q_1, \dots, Q_n\}$  s.t.  $Q_i \subseteq U$ , and an integer  $k \leq |Q|$ , a set  $I \subseteq Q$  is a  $k$ -intersecting set if:  $|I| = k, \bigcap_{Q \in I} Q \subseteq M$ , and  $\bigcap_{Q \in I} Q \neq \emptyset$ .

**Theorem: K-SET-INTERSECTION is NP-complete (reduction from 3-SAT).**

**Algorithm Rationale**

Let  $T_m = \{R_{m,i} | R_{m,i} = (U \setminus M) \setminus (Q_i \setminus M) \text{ and } m \in M, Q_i\}$

If  $\exists R_{m,1} \cup \dots \cup R_{m,k} = U \setminus M$  then  $\overline{R_{m,1}} \cap \dots \cap \overline{R_{m,k}} = \emptyset$

Since  $\overline{R_{m,i}} = (Q_i \setminus M)$  and  $m \in Q_i$  then  $m \in Q_1 \cap \dots \cap Q_k$  and  $Q_1 \cap \dots \cap Q_k \subseteq M$

**k-Set-Intersection Approximation**

- Greedy algorithm
  - Uses Set Cover greedy approximation algorithm at its core
- Given  $(U, M, Q, k)$  do:
  - Step 1:  $\forall m \in M, T_m = \{(U \setminus M) \setminus (Q_i \setminus M) : m \in Q_i\}$
  - Step 2: Run k-SET-COVER greedy algorithm on  $(U \setminus M, T_m, k)$ 
    - 2a: Pick  $R \in T_m$  with the maximum uncovered elements
    - 2b: Take the union of every set picked in 2a
    - 2c: If the union is  $U \setminus M$  go to step 3, else if we picked less than  $k$  sets go to 2a, else repeat for another  $m \in M$
  - Step 3: For every set  $(U \setminus M) \setminus (Q_i \setminus M)$  in the set cover, add  $Q_i$  in the intersecting set

**Approximation Algorithm: APRX-SFW**

- Adopt k-Set-Intersection Approximation:
  - $U = S$  the set of servers
  - $Q = \{Q_1, \dots, Q_k\}, Q_i \subseteq S$  is the quorum system
  - $M \subseteq S$  the servers that replied with the latest value
  - $k$  the number of quorums required by the predicates
- Log-Approximation
  - Invalidates RP and WP a factor of  $\log|S|$  times
- What does it mean for SFW?
  - Extra Communication Rounds (esp. for writes)
  - Slower acceptance of a new value
  - Does not affect correctness

## Algorithm CWFR

**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  $Q$
- $QView_c(1)$  - Fast: return maxTag of current iteration
- $QView_c(2)$  - remove servers with maxTag and re-evaluate
- $QView_c(3)$  - Slow: propagate and return maxTag

**Server Protocol: passive role**

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

**Quorum Views in SWMR**

If a read operation receives replies from quorum  $Q$  will observe one of the following:

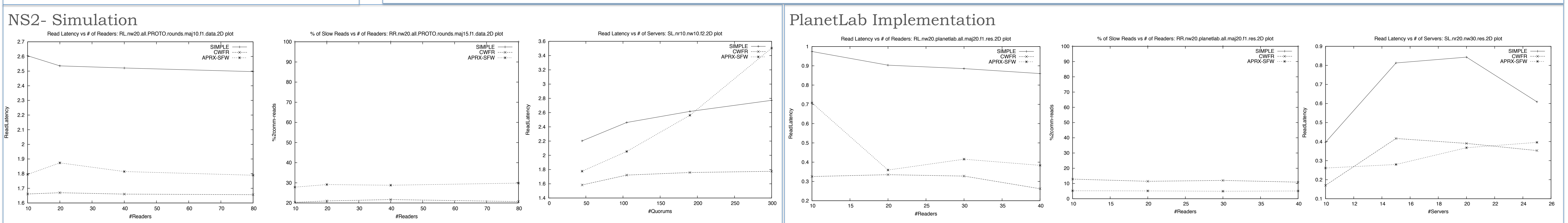
- Full quorum contains  $\max(\text{tag})$  => Write potentially completed
- Every quorum contains an old tag => Write of  $\max(\text{tag})$  not completed
- Every server of an intersection replied with  $\max(\text{tag})$  => either write not completed c(i), or write completed in  $Q_c$  c(ii)

a,b = write status Determinable => FAST  
c = write status Undeterminable => SLOW

**Uncover the past in MWMR**

Discover the latest potentially completed write For tags  $\langle 1, \rangle, \langle 2, \rangle, \langle 3, \rangle$ :

- $\langle 3, \rangle$  not completed (servers possibly contained  $\langle 2, \rangle, \langle 1, \rangle$ )
- $\langle 2, \rangle$  not completed (servers possibly contained  $\langle 1, \rangle$ )
- $\langle 1, \rangle$  potentially completed



ΚΥΠΡΙΑΚΗ ΔΗΜΟΚΡΑΤΙΑ

ΕΥΡΩΠΑΪΚΗ ΕΝΩΣΗ

Η ΔΕΣΜΗ 2009-10 ΣΥΓΧΡΗΜΑΤΟΔΟΤΕΙΤΑΙ ΑΠΟ ΤΗΝ ΚΥΠΡΙΑΚΗ ΔΗΜΟΚΡΑΤΙΑ ΚΑΙ ΤΟ ΕΥΡΩΠΑΪΚΟ ΤΑΜΕΙΟ ΠΕΡΙΦΕΡΕΙΑΚΗΣ ΑΝΑΠΤΥΞΗΣ ΤΗΣ ΕΕ

UCONN