







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

Seeking Fast Operations in MWMR Atomic Register Implementations

Nicolas Nicolaou

University of Cyprus & University of Connecticut

Funded by the Cyprus Research Promotion Foundation and co-funded by the Republic of Cyprus and the European Regional Development Fund

What is a Distributed Storage System?



Complexity Measure-Operation Latency

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



communication nounds (round-trips)

What was known...



The Era of Fast Implementations...



Exploring the possibilities in the MWMR

Can we obtain fast write and read operations in the MWMR model?

Achievements...

- Introduced new technique for value ordering
 - Server Side Ordering (SSO)
- Devised an Atomic Register implementation in the MWMR model
 - Enables Fast Writes and Reads -- first such algorithm
 - Allows Unbounded Participation

SFW Algorithm – An intuitive idea

- Value versioning
 - Associate a tag with each value
- Traditional Approaches
 - Increment tag (version) at the writer
- SSO Approach
 - Increment tag (version) at the server
 - Avoid discovery of the latest version from the writer

Traditional Writer-Server Interaction



SSO Writer-Server Interaction



Definition: n-wise Quorum Systems

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

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



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 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

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

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

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 $\mathbb{Q} = \{Q_1, \ldots, Q_n\}$ s.t. $Q_i \subseteq U$, and an integer $k \leq |\mathbb{Q}|$, a set $I \subseteq \mathbb{Q}$ 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).

The Good News...

Approximation Algorithm (APRX-SFW)

- Greedy Algorithm
- Log-approximation
 - log|S| the number of slow operations
- Algorithm CWFR
 - Based on Quorum Views
 - SWMR prediction tools
 - Fast operations in General Quorum Systems
 - Trades Speed of Write operations
 - Two Round Writes

