

institute i dea networks









NAZARBAYEV UNIVERSITY

Making "Fast" Atomic Operations Computationally Tractable

Antonio Fernández Anta, Nicolas Nicolaou, and Alexandru Popa

Developing the Science of Networks

Problem Statement



Implementing a fault-tolerant shared object in an asynchronous, message-passing environment:

- Availability + Survivability => use redundancy
- Asynchrony + Redundancy => concurrent operations
- Behavior of concurrent operations => consistency semantics
 - Safety, Regularity, Atomicity [Lamport86]

System Model: Definitions

Components	 Clients: 1 writer & R readers (SWMR) Servers: S replica hosts 		
Operations	 write(v): updates the object value to v read(): retrieves the object value Well-Formedness (only a single operation at a time) 		
Communication	 Asynchronous Message-Passing Reliable Channels (messages are not lost or altered) 		
Failures	 Crashes Any reader or the writer Up to minority of servers 		

3

91/61/1

Nicolas Nicolaou



Consistency Model: Atomicity/Linearizability

- Provides the illusion that operations happen in a sequential order
 - a read returns the value of the preceding write
 - a read returns a value at least as recent as that returned by any preceding read

Complexity Measure



Nicolas Nicolaou

91/61/1

5

- [Attiya, Bar-Noy, Dolev 1996] (Dijkstra Prize 2011)
- Order Operations by using <ts, v> pairs.

Reader Protocol (2 phases)

- Phase 1:
 - Send read to all
 - Collect <ts,v> from a majority
 - Discover max(<ts,v>)
- Phase 2:
 - Send max(<ts,v>) to all
 - Collect ack from a majority and return v

Reads must Write! (2 round-trips)

Writer Protocol

- ts++ //increment ts
- Send <ts,v> to all
- Wait for a majority to reply

Server Protocol (Upon rcv a msg m from i)

- if m.ts > ts
 - Update local <ts,v>
- Send local <ts,v> to i

Nicolas Nicolaou

Algorithm FAST: Reducing communication

- [Dutta et al. 2004, 2010] made a nice observation
- Consider f=1, S=5 and let operation to communicate with S-f servers instead of majority





Algorithm FAST: Reducing communication

- [Dutta et al. 2004, 2010] made a nice observation
- Consider f=1, S=5 and let operation to communicate with S-f servers instead of majority



After 1 read

- At least S-2f servers
- Each replies to both {w, r1}

8



Algorithm FAST: Reducing communication

- [Dutta et al. 2004, 2010] made a nice observation
- Consider f=1, S=5 and let operation to communicate with S-f servers instead of majority



Nicolas Nicolaou

9





Algorithm FAST: Single round-trip operations

• Constructed a predicate that allowed all writes and reads to complete in a single round-trip

 $\exists \alpha \in [1, R+1] \land MS \subseteq S \text{ s.t.}$

$$\forall s \in MS, s.ts = maxTS \land |MS| \ge S - \alpha f \land |\bigcap_{s \in MS} s.seen| \ge \alpha$$

Server Protocol (Upon rcv a msg m from i)

- if m.ts > ts
 - Update local info
 - Reset seen set to {i}
- else
 - Add i in seen set
- Send ts and seen set to i

Writer Protocol (1 phase)

(Same as ABD)

Reader Protocol (1 phase)

- Send read to all
- Collect <ts,v> and seen sets from S-f
- Discover maxTS = max(ts)
- If predicate is true:
 - return maxTS
- else
 - return maxTS-1



• On the number of readers:

$$R < \frac{S}{f} - 2$$

- On the number of writers
 - Impossible in the MWMR model

• What about the computation?

11





How hard is it to compute the predicate?

Problem 1 (Predicate Formalization)

Input

Two sets

 $U_1 = \{s_1, \dots, s_n\}, U_2 = \{p_1, \dots, p_k\} \text{ s.t. } \forall s_i \in U_1, s_i \subseteq U_2$

• Two integers α, f s.t. $n - \alpha f \ge 1$

Output

Is there a set

$$M \subseteq U_1$$
 s.t. $\left| \bigcap_{s \in M} s \right| \ge \alpha$ and $|M| > n - \alpha f$

Problem equivalence to the predicate

$$\exists \alpha \in [1, R] \land MS \subseteq S \quad \text{s.t.}$$

$$\forall s \in MS, s.ts = maxTS \land |MS| \ge S - \alpha f \land |\bigcap_{s \in MS} s.seen| \ge \alpha$$

- U₁ : the set of all the seen sets collected
- U₂ : set of reader and the writer identifiers
- M : MS in the predicate

13

91/61/1

Vicolas Nicolaou

• α , f: the respective α and f in the predicate

Maximum Biclique Problem (MBP)

s1

s2

s3

pl

p2

р3

Definition MBP:

Input

- A bipartite graph G = (X, Y, E)
- A positive integer c

Output

 $A = \{s1, s3\}, B = \{p1, p2\}$ c = 4 Are there two sets $A \subseteq X, B \subseteq Y$ s.t. $\forall a \in A, \forall b \in B, (a, b) \in E \text{ and } |E| = |A| * |B| \ge c$

MBP is NP-complete [Peeters 01]

14

Problem 1 is NP-Hard - Reduction Idea

Input Transformation

- Given the graph G = (X, Y, E)
 - Set $U_1 = X$
 - Set $U_2 = Y$
 - $(s_i, p_j) \in E \iff p_j \in s_i$
- Set $c = \alpha(n \alpha f)$
- Problem 1 returns true if
 - exists $M \subseteq U_1$ and $|M| = n \alpha f$

 $\exists P \subseteq U_2, |P| = \alpha \text{ s.t. } \forall s \in M, \forall p \in P, (s, p) \in E$

• In this case $c = |M| * |P| = \alpha(n - \alpha f)$ and MBP is true



Can we overcome NP-Hardness?

Observation

To avoid the excessive computation we need to avoid the set manipulation.

Question

Can we preserve atomicity if we know how many and not which processes read the latest value?

16





Writer Protocol (1 phase) (Same as ABD)

Nicolas Nicolaou

1/19/16

imtitute networks

 Consider f=1, S=5 and let operation to communicate with S-f servers



18

1/19/16

Nicolas Nicolaou

 Consider f=1, S=5 and let operation to communicate with S-f servers



After 1 read

- At least S-2f servers
- Each replies to 2 processes

19

- [Dutta et al. 2004] made a nice observation
- Consider f=1, S=5 and let operation to communicate with S-f servers instead of majority



After 1 read

- At least S-2f servers
- Each replies to 2 processes

After 2 reads

- At least S-3f servers
- Each replies to **3** processes

Nicolas Nicolaou

20

After k-1 reads

- At least S-kf servers
- Each replies to k processes
- How big can k be?

$$S - kf > f \Rightarrow k < \frac{S}{f} - 1$$

• Since k is the number of processes then

$$R+1 < \frac{S}{f} - 1 => R < \frac{S}{f} - 2$$

21





Computing the new predicate in linear time



Linear Complexity: O(S)

Algorithm

- •Given b[1...R+1]
- •For each s that replied
 - if s.ts = maxTS
 - Increment b[s.views]
- •For bucket α in R+1 to 2
 - $\text{ If } b[\alpha] \ge S \alpha f \text{ return true}$
 - Else "empty" $b[\alpha]$ in $b[\alpha-1]$
 - b[α-1] += b[α]



Complexity Comparison

- Server Messages
 - ABD: [<ts, v>]
 - FAST: [<ts,v>, set seen]
 - ccFAST: [<ts,v>, int views]

Algorithm	WR	$\mathbf{R}\mathbf{R}$	WC	RC	WB	RB
ABD	1	2	O(1)	$O(\mathcal{S})$	$O(\lg V)$	$O(\lg V)$
Fast	1	1	O(1)	$O(\mathcal{S} ^2 \cdot 2^{ \mathcal{S} })$	$O(\lg V)$	$\Theta(\mathcal{S} + \lg V)$
CCFAST	1	1	O(1)	$O(\mathcal{S})$	$O(\lg V)$	$O(\lg \mathcal{S} + \lg V)$

WR/RR: write/read round-trips WC/RC: write/read computation demands WB/RB: write/read message size in bits

Conclusions

- We showed that FAST is not computationally tractable
- Proposed a new predicate that
 - Preserves "fast" behavior from the operations
 - Reduces the messages sizes
 - Can be computed in linear time
 - Preserves Atomicity
- Presented an algorithm that computes the proposed predicate in linear time

Redefined "fastness" for Atomic Operations

Nicolas Nicolaou

Thank you !



25

91/61/1

Nicolas Nicolaou

institute

62

Reduction Example

• Consider Previous Example (f=1, S=5)

26

91/61/1

Nicolas Nicolaou





Consistency Model: Atomicity/Linearizability

- Provides the illusion that operations happen in a sequential order
 - a read returns the value of the preceding write
 - a read returns a value at least as recent as that returned by any preceding read



- [Attiya, Bar-Noy, Dolev 1996] (Dijkstra Prize 2011)
- Order Operations by using <ts, v> pairs.





Writer Protocol

- ts++ //increment ts
- Send <ts,v> to all
- Wait for a majority to reply

Server Protocol (Upon rcv a msg m from i)

- if m.ts > ts
 - Update local <ts,v>
- Send local <ts,v> to i

28

- [Attiya, Bar-Noy, Dolev 1996] (Dijkstra Prize 2011)
- Order Operations by using <ts, v> pairs.



29

1/19/16

Vicolas Nicolaou

Writer Protocol

- ts++ //increment ts
- Send <ts,v> to all
- Wait for a majority to reply

- [Attiya, Bar-Noy, Dolev 1996] (Dijkstra Prize 2011)
- Order Operations by using <ts, v> pairs.



Writer Protocol

- ts++ //increment ts
- Send <ts,v> to all
- Wait for a majority to reply

Nicolas Nicolaou

30

- [Attiya, Bar-Noy, Dolev 1996] (Dijkstra Prize 2011)
- Order Operations by using <ts, v> pairs.



31

1/19/16

Nicolas Nicolaou

Reader Protocol (2 phases)

- Phase 1:
 - Send read to all
 - Collect <ts,v> from a majority
 - Discover max(<ts,v>)
- Phase 2:
 - Send max(<ts,v>) to all
 - Collect ack from a majority and return v

Reads must Write!

Consistency Semantics [Lamport86]



institute iMde network

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 communicaiton