



institute
imdea
networks

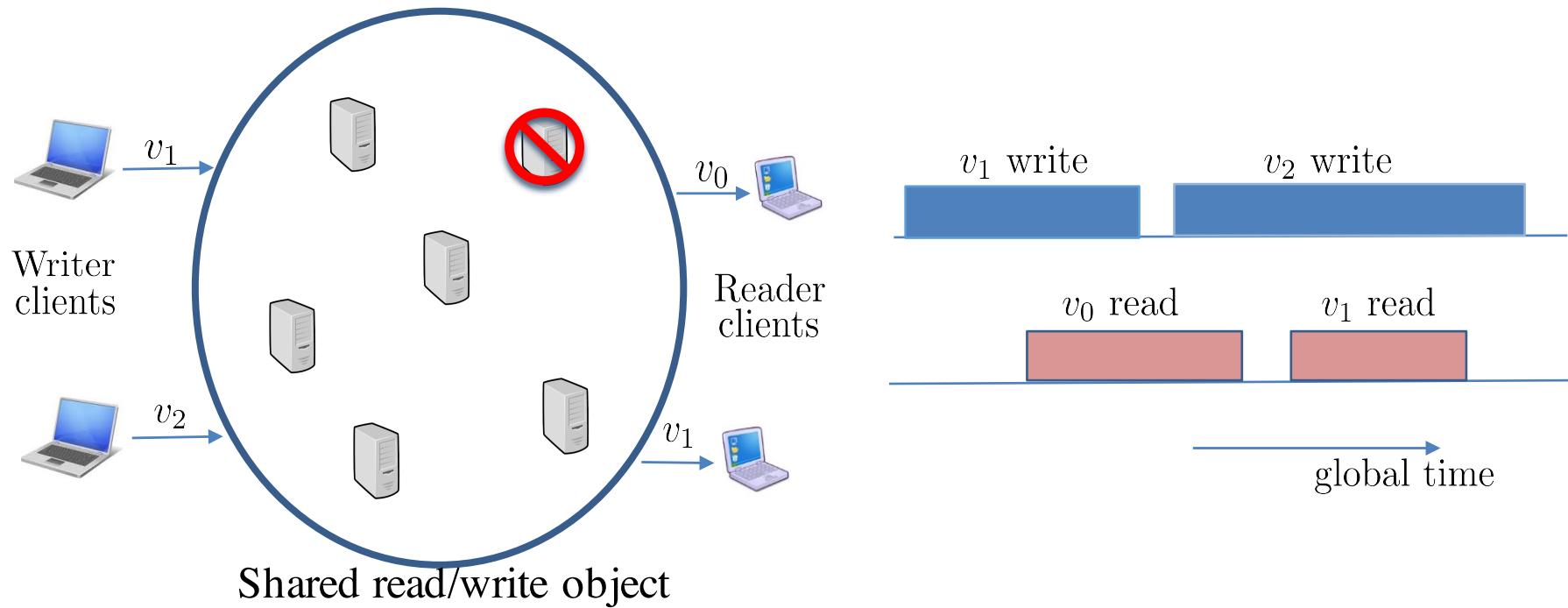


CoVer-ability: Consistent Versioning in Asynchronous, Fail-Prone, Message-Passing Environments

Nicolas Nicolaou, Antonio Fernández Anta, and Chryssis Georgiou

[Developing the
Science of Networks]

Distributed Read/Write Objects

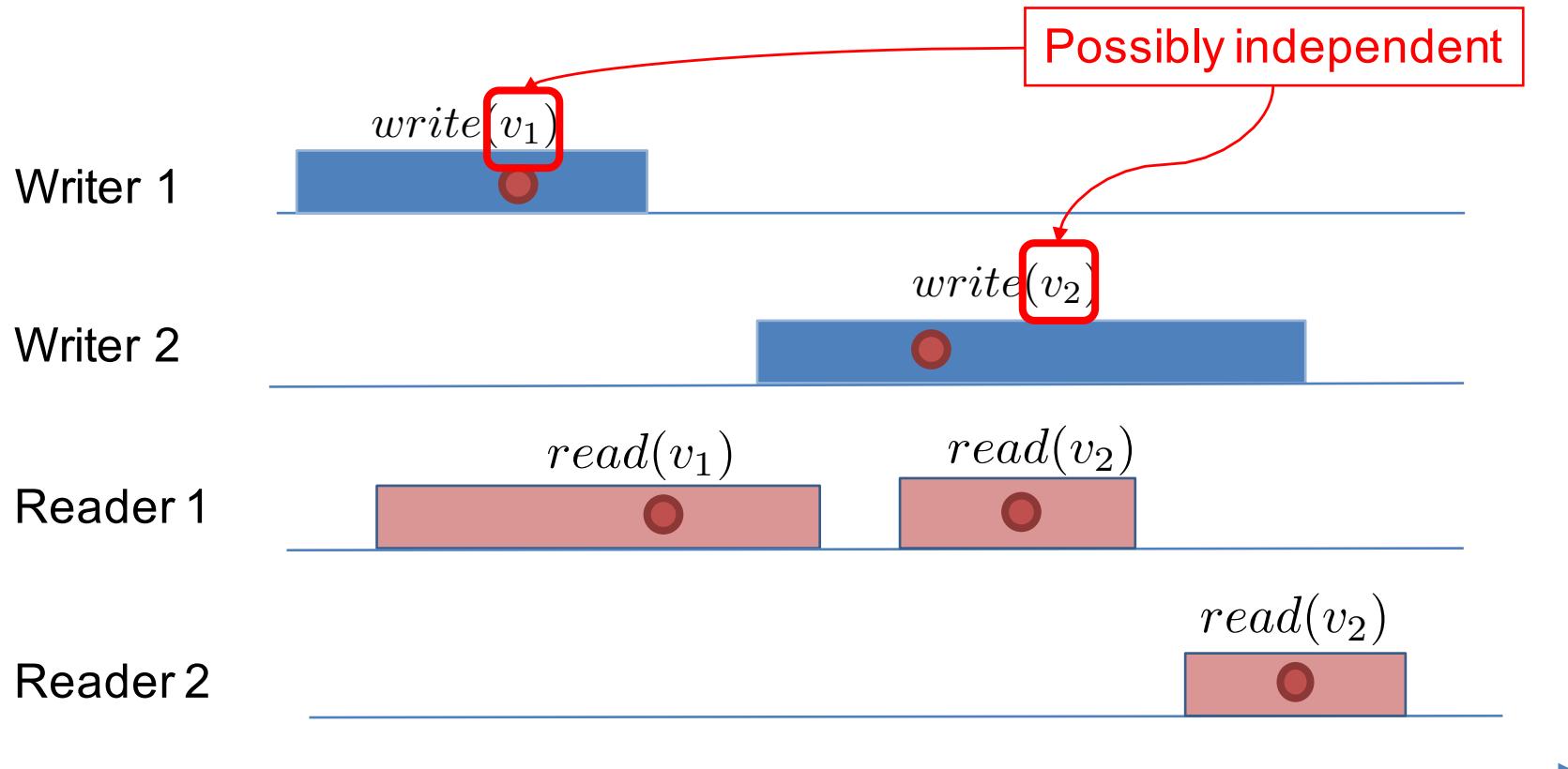


Implementing a **fault-tolerant** shared object in an **asynchronous, message-passing** environment:

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

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**



Consistency on Object Evolution

- Particular objects may need **consistent state evolution**

- File objects
 - Distributed databases
 - Bulleting boards



- Expected Behavior: If a write w_1 completes and writes version i , then any write w_2 that succeeds w_1 must write a version j s.t.
 - version j obtained by modifying version i
 - version j obtained by modifying version k and there is a sequence of modifications that lead from i to k .



Not preserved with existing consistency semantics

How can we achieve such guarantees?

- Standard practices use:
 - Compare and Swap (CAS)
 - Linked-Load / Store-Conditional (LL/SC)
- Such primitives can solve consensus
 - Impossible in the asynchronous, message-passing, fail prone environment [FLP85]

Question:

Can we provide **versioning guarantees** in an asynchronous, message-passing, fail-prone environment **using weaker primitives, like read/write registers?**

System Model: Definitions

Components

- n processes
- Asynchronous Computation

Communication

- Asynchronous
- Message-Passing
- Reliable Channels (messages are not lost or altered)

Failures

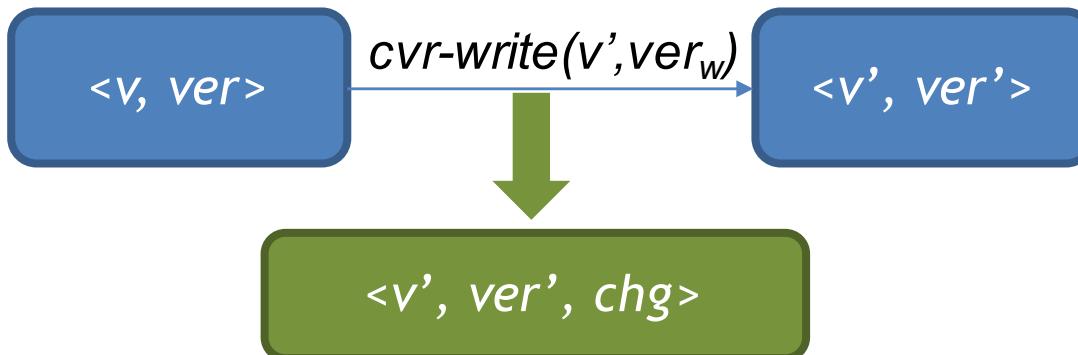
- Crashes

Object Type: Versioned Read/Write Register

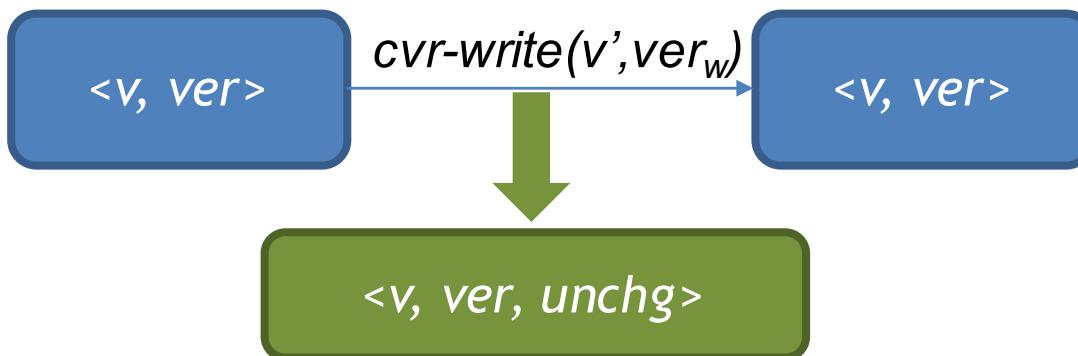
- *Versions*: a totally ordered set of versions
- Domain of Values: *Values x Versions*
 - Each value is assigned with a version
- R/W Register Operations
 - `write(v)`: updates the object value to v
 - `read()`: retrieves the object value
- Versioned R/W Register Operations
 - `cvr-write(v,ver)`: attempt to change the value of the object **with version** ver and **returns whether successful or not**
 - `cvr-read()`: retrieves the object value **and version**

Versioned Register Operations

- Let $\langle val, ver \rangle$ be the state of the register
- Write operations on the versioned register are:
 - Successful write given that $ver_w = ver$:

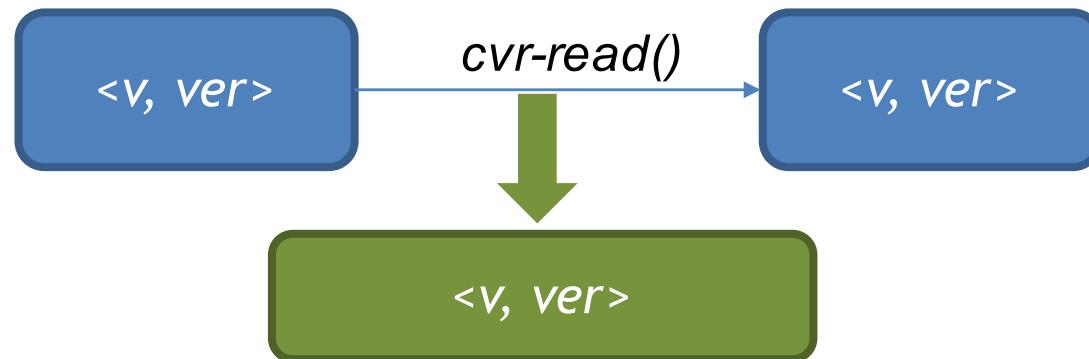


- Unsuccessful write given that $ver_w \neq ver$:



Versioned Register Operations

- Read operations on the versioned register:



- Notation we use:
 - Successful write: $cvr-w(ver)[ver']$
 - Unsuccessful write: $cvr-w(ver)[ver',unchg]$

Definition: Execution Validity

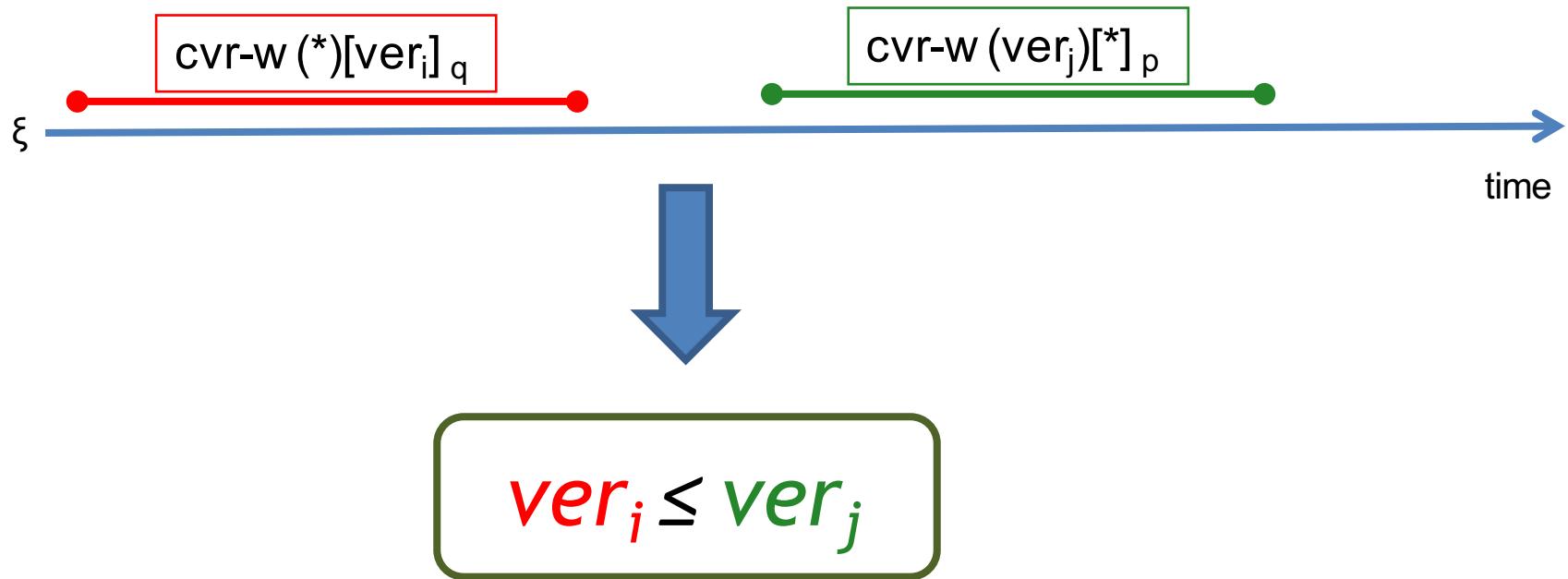
- Version increment
 - For any successful write $\text{cvr-w}(\text{ver})[\text{ver}']$, $\text{ver}' > \text{ver}$
- Version uniqueness
 - For any two successful writes $\text{crv-w}(\text{*})[\text{ver}']$ and $\text{crv-w}(\text{*})[\text{ver}'']$, then $\text{ver}' \neq \text{ver}''$
- Version connection
 - For each successful write $\text{cvr-w}(\text{ver}_{k-1})[\text{ver}_k]$, there is a sequence of versions $\text{ver}_0, \text{ver}_1, \dots, \text{ver}_k$ s.t. there is a successful write $\text{cvr-w}(\text{ver}_i)[\text{ver}_{i+1}]$, for $0 < i < k-1$

Definition: Coverability

- An **execution** is **coverable** with respect to some order $<_\xi$ if it is **valid** and it satisfies:
 - Version Consolidation
 - Version Continuity
 - Version Evolution

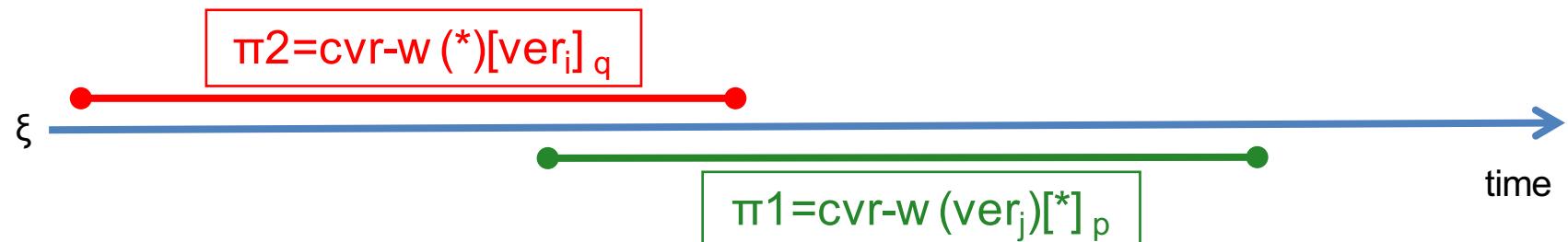
Consolidation in Execution ξ

- If two successful writes $\pi_1 = \text{crv-w}(*)[\text{ver}_i]$ and $\pi_2 = \text{crv-w}(\text{ver}_j)[*]$, s.t. $\pi_1 \rightarrow \pi_2$ in ξ , then $\text{ver}_i \leq \text{ver}_j$ and $\pi_1 <_{\xi} \pi_2$



Continuity in Execution ξ

- If $\pi_1 = \text{cvr-w}(\text{ver}_j)[*]$ is a successful write, then there exists successful write $\pi_2 = \text{cvr-w}[*][\text{ver}_i]$, s.t. $\pi_2 <_{\xi} \pi_1$ and $\text{ver}_i = \text{ver}_j$, or $\text{ver}_j = \text{ver}_0$



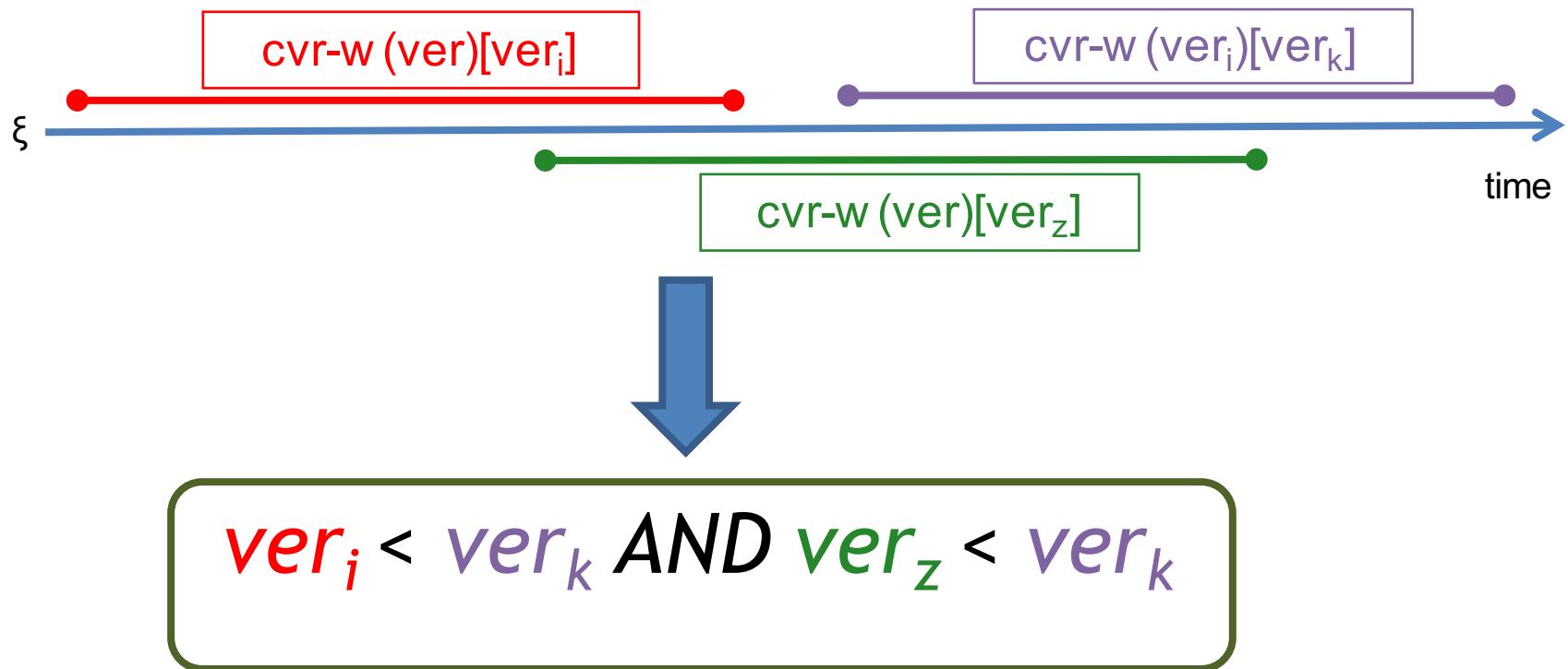
$\pi_2 <_{\xi} \pi_1$ AND $\text{ver}_i = \text{ver}_j$

OR

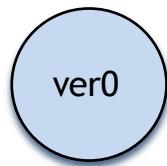
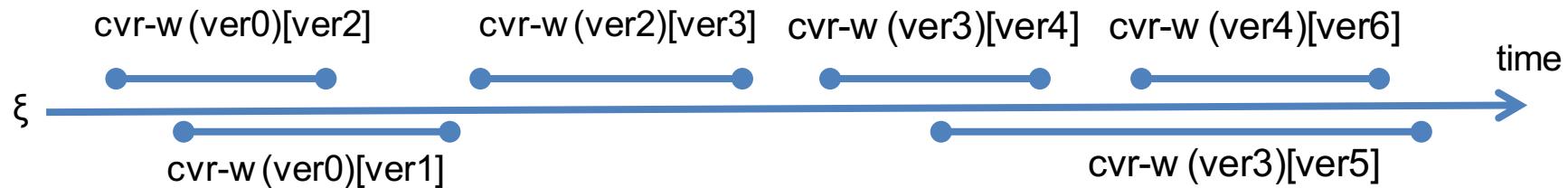
$\text{ver}_j = \text{ver}_0$

Evolution in Execution ξ

- If $\pi_1 = \text{crv-w}(\text{ver}_i)[\text{ver}_{i+1}]$ and $\pi_2 = \text{crv-w}(\text{ver}_j)[\text{ver}_{j+1}]$, two successful writes, s.t. $\text{ver}_i \leq \text{ver}_j$, then $\text{ver}_{i+1} < \text{ver}_{j+1}$



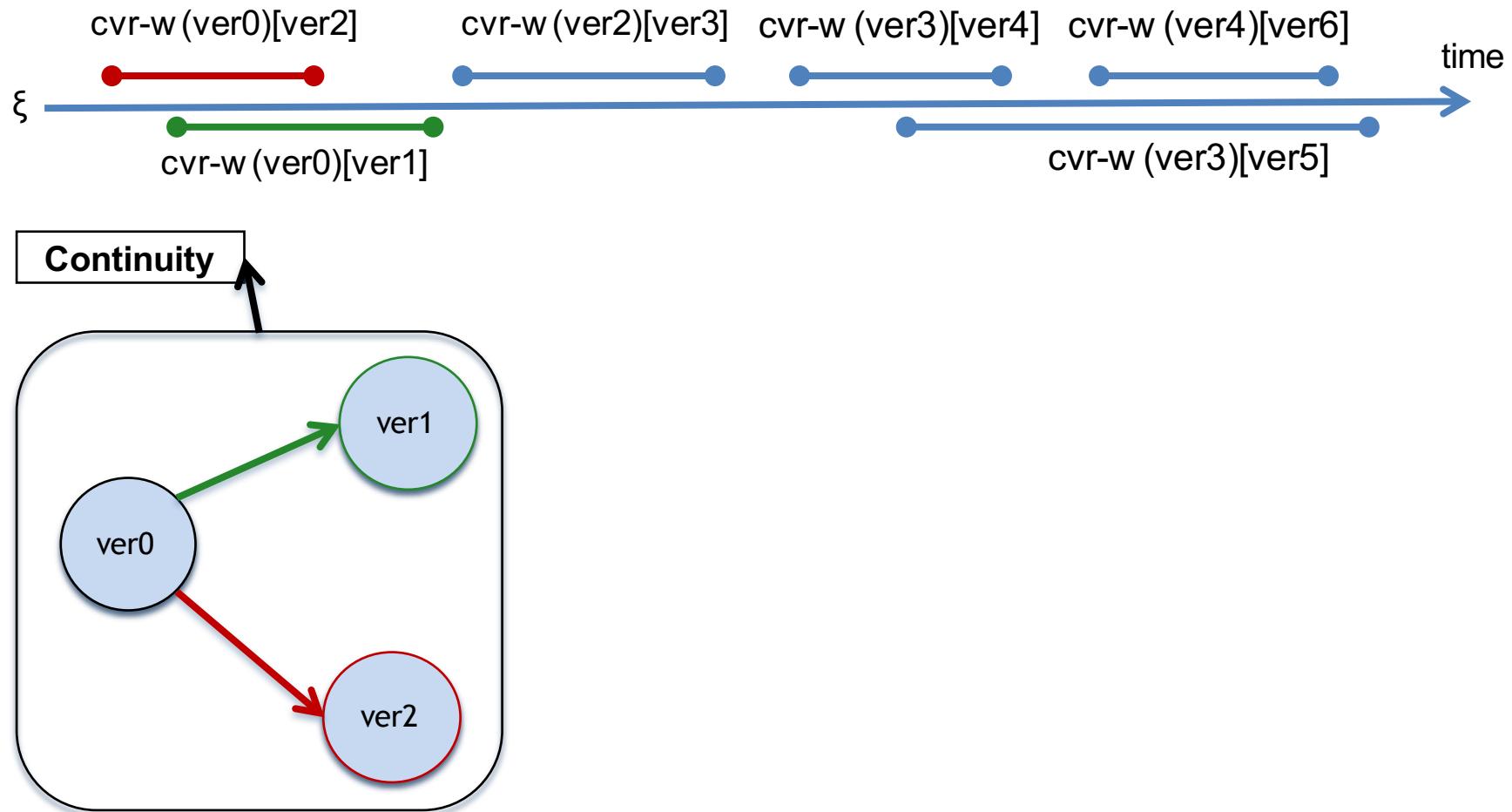
Coverable Execution Tree



Total Order:

$\text{cvr-w}(\text{ver0})[\text{ver1}] <_{\xi} \text{cvr-w}(\text{ver0})[\text{ver2}] <_{\xi} \text{cvr-w}(\text{ver2})[\text{ver3}] <_{\xi} \text{cvr-w}(\text{ver3})[\text{ver5}] <_{\xi} \text{cvr-w}(\text{ver3})[\text{ver4}] <_{\xi} \text{cvr-w}(\text{ver4})[\text{ver6}]$

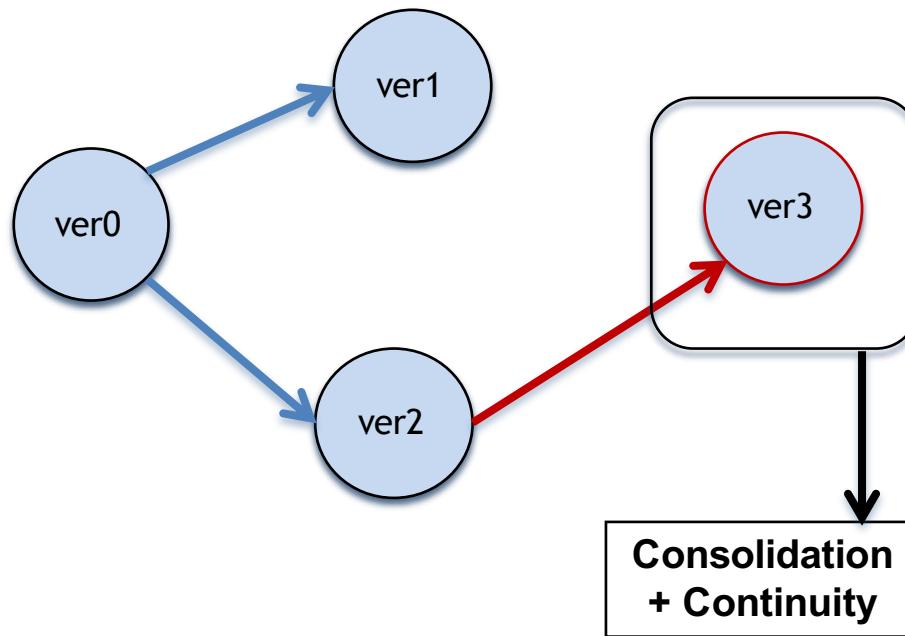
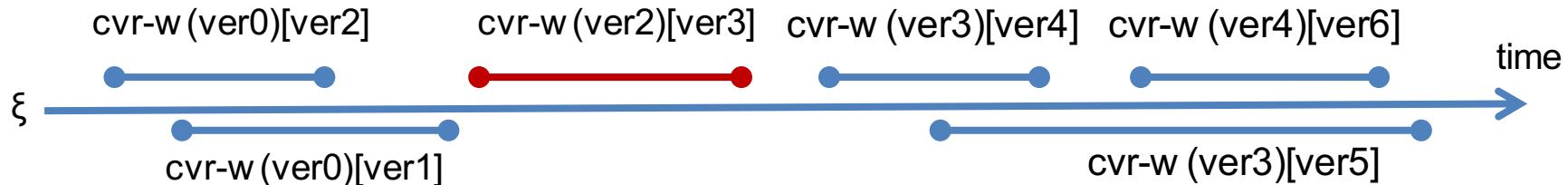
Coverable Execution Tree



Total Order:

$cvr-w(\text{ver0})[\text{ver1}] <_{\xi} cvr-w(\text{ver0})[\text{ver2}] <_{\xi} cvr-w(\text{ver2})[\text{ver3}] <_{\xi} cvr-w(\text{ver3})[\text{ver5}] <_{\xi} cvr-w(\text{ver3})[\text{ver4}] <_{\xi} cvr-w(\text{ver4})[\text{ver6}]$

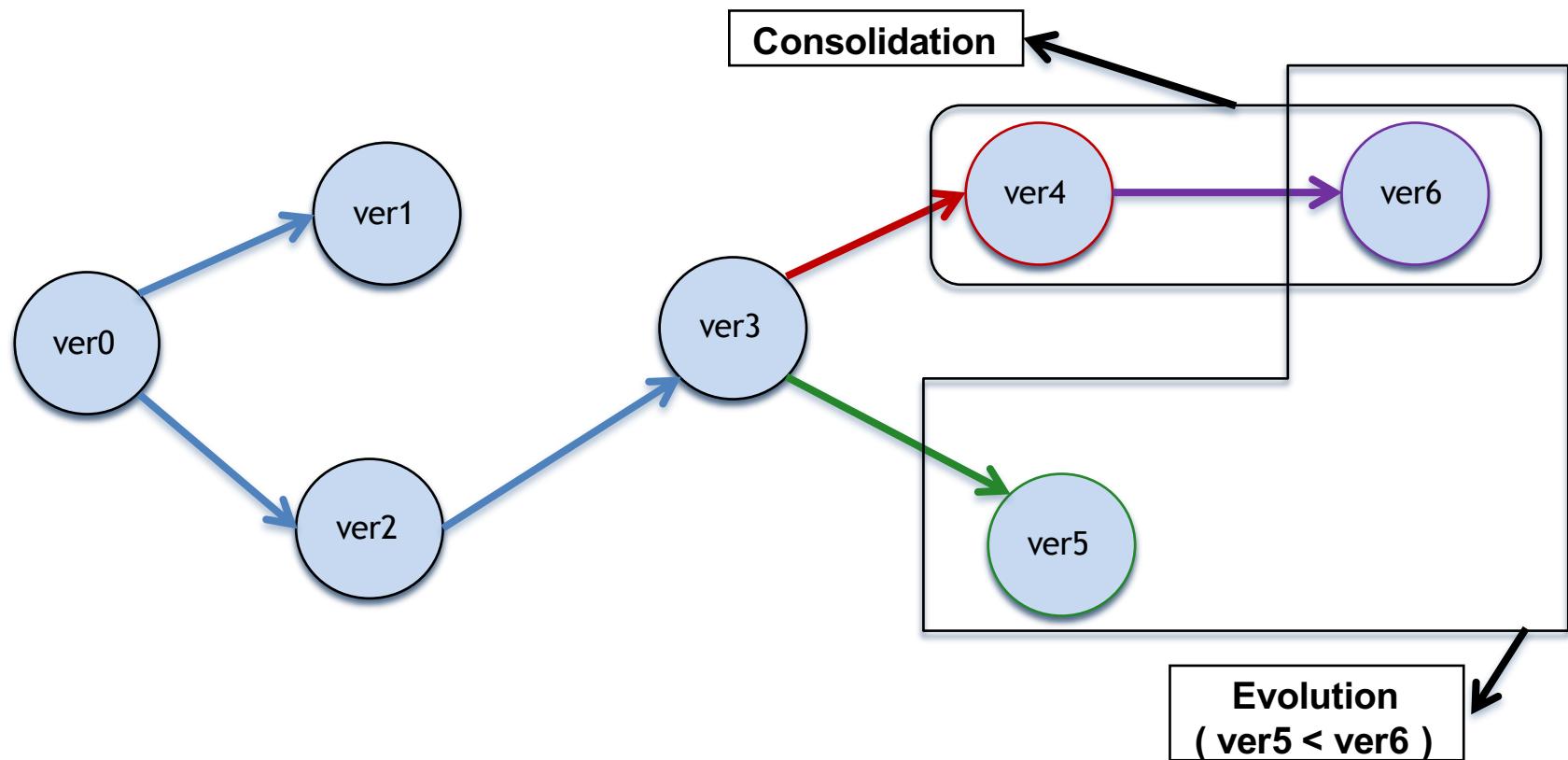
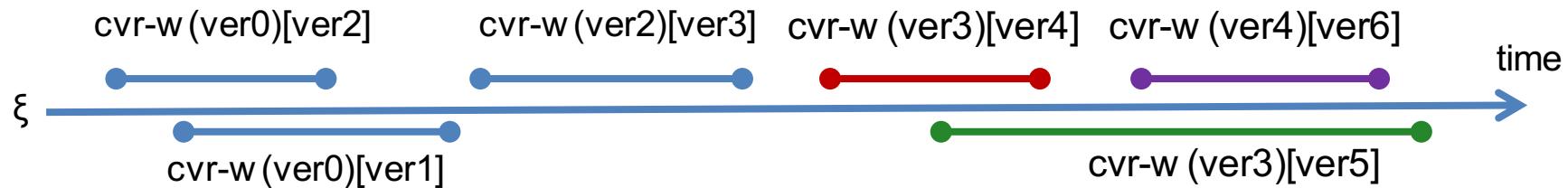
Coverable Execution Tree



Total Order:

$\text{cvr-w (ver0)[ver1]} <_{\xi} \text{cvr-w (ver0)[ver2]} <_{\xi} \text{cvr-w (ver2)[ver3]} <_{\xi} \text{cvr-w (ver3)[ver5]} <_{\xi} \text{cvr-w (ver3)[ver4]} <_{\xi} \text{cvr-w (ver4)[ver6]}$

Coverable Execution Tree



Total Order:

$\text{cvr-w}(\text{ver0})[\text{ver1}] \lessdot_{\xi} \text{cvr-w}(\text{ver0})[\text{ver2}] \lessdot_{\xi} \text{cvr-w}(\text{ver2})[\text{ver3}] \lessdot_{\xi} \text{cvr-w}(\text{ver3})[\text{ver5}] \lessdot_{\xi} \text{cvr-w}(\text{ver3})[\text{ver4}] \lessdot_{\xi} \text{cvr-w}(\text{ver4})[\text{ver6}]$

Definition: Coverable Atomic Register

- A versioned register is coverable and atomic , referred as **coverable atomic register** , if any execution on the register satisfies:
 - **Atomicity**, and
 - **Coverability** with respect to the total order imposed by atomicity on the successful writes.

Coverable vs Ranked Registers

Ranked Registers:

- Registers that support operations where value is associated with a rank
- Writes may commit or abort
- Writes abort => exists a rank higher than the one to be written
- May write any arbitrarily large rank

Theorem:

There is no algorithm that implements a coverable register using a set of ranked registers.

Implementing Coverable Atomic R/W Registers

- Enhance ABD algorithm [ABD96]
 - Register replicated to S servers
 - Use $\langle wid, ts, v \rangle$ triples to order operations
 - $tag = \langle wid, ts \rangle$ is used as the version of the object
 - tags are compared lexicographically

Server Protocol (Upon rcv a msg m from i)

- if $m.type \neq \text{Query} \ \&\& \ m.tag > tag$
 - Local $\langle tag, v \rangle = \langle m.tag, m.v \rangle$
 - Send local $\langle tag, v \rangle$ to i

Implementing Coverable Atomic R/W Registers

cvr-write(val,ver) at w_i

- Phase 1:
 - Send **Query** to all
 - Collect $\langle tag, v \rangle$ from a majority
 - Discover $\max(\langle tag, v \rangle)$
- Phase 2:
 - If ($ver == \max(tag)$)
 - $\langle tag', v' \rangle = \langle w_i, \max(tag.ts)+1, val \rangle$
 - $flag = chg$
 - Else
 - $\langle tag', v' \rangle = \max(\langle tag, v \rangle)$
 - $flag = unchg$
 - Send **Write** $\max(\langle tag', v' \rangle)$ to all
 - Collect ack from a majority and return $(\langle tag', v' \rangle, flag)$

- Read is similar to the write operation except that in **Phase 2** it propagates the $\max(\langle tag, v \rangle)$ discovered.
- Read **returns the $\max(\langle tag, v \rangle)$** discovered during Phase 1

Application: Weak RMW

- Implement **weaker RMW semantics** using coverable R/W registers
 - Not all RMW operations may modify the value of the shared RMW object
 - In the case of concurrency at least one will RMW the object

weakRMW(F)

- $\langle \text{oldval}, \text{lcver} \rangle = \text{cvr-read}()$
- $\text{newval} = F(\text{oldval})$
- $\langle \text{lcval}, \text{lcver}, \text{flag} \rangle = \text{cvr-write}(\text{newval}, \text{lcver})$
- If ($\text{flag} == \text{chg}$)
 - return $\langle \text{lcval}, \text{success} \rangle$
- Else
 - return $\langle \text{lcval}, \text{failure} \rangle$

Application: Concurrent File Objects

- File operations: `get()`, `revise()`
- A `get()` operation **returns the latest version** of the file
- A `revise()` operation **specifies the version** of the file to be modified
 - Prevent revising older versions

Get()

- `<lcval,lcver> = cvr-read()`
- `return <lcval, lcver>`

Revise(v, ver)

- `<lcval, lcver, flag> = cvr-write(v, ver)`
- If (`flag == chg`)
 - `return OK`
- `return <lcval, lcver>`

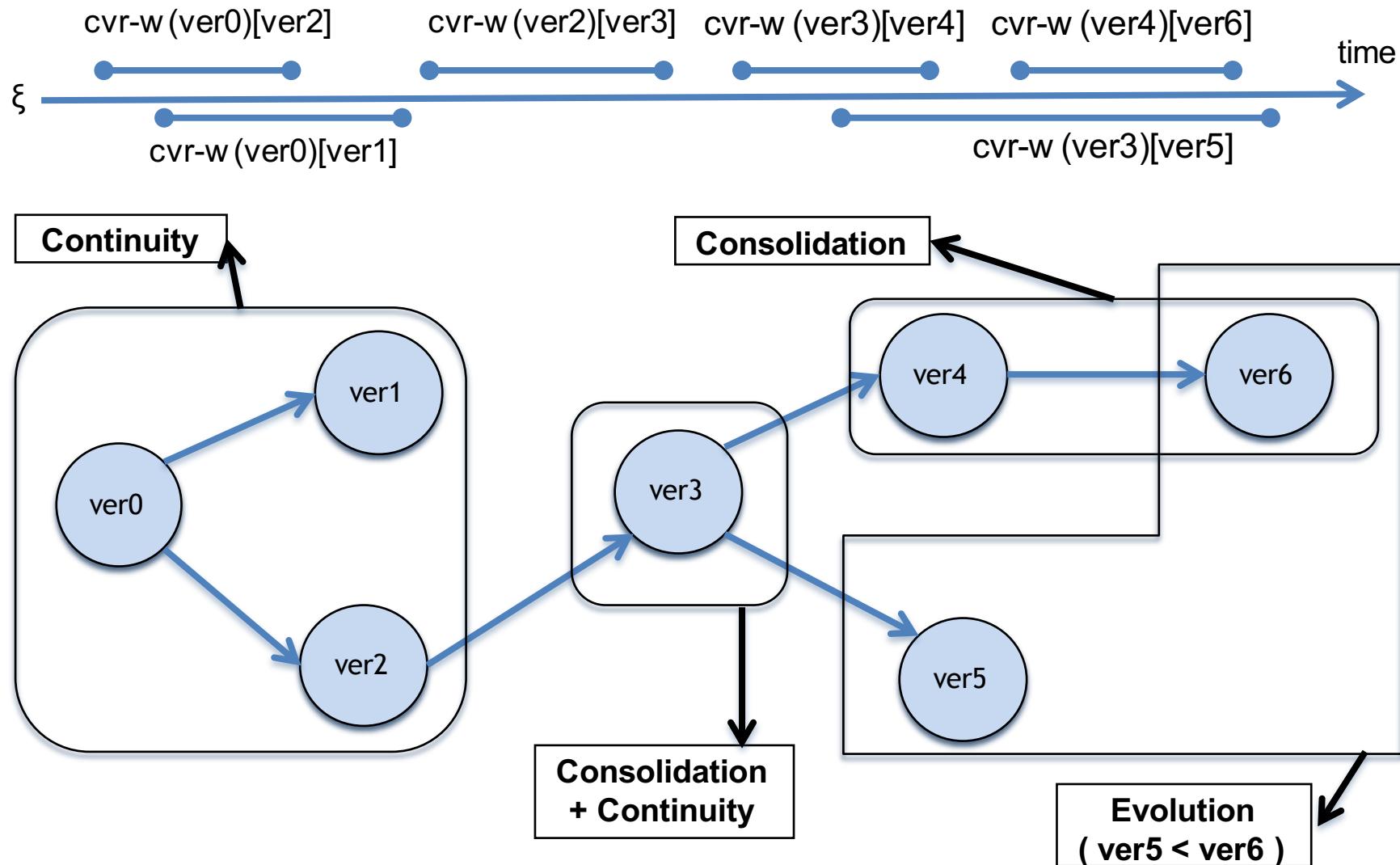
Conclusions

- We defined a new object type: **versioned object**
- We defined **Coverability** – a **new property** for concurrent versioned objects
- We defined **Atomic Coverable R/W Registers**
- We compared ACR similar objects (Ranked Registers)
- Demonstrated an **implementation** of ACR
- We presented two applications
 - Weak RMW objects
 - Distributed File objects

Thank you !



Coverable Execution Tree



Total Order:

$cvr-w(\text{ver0})[\text{ver1}] \leq_{\xi} cvr-w(\text{ver0})[\text{ver2}] \leq_{\xi} cvr-w(\text{ver2})[\text{ver3}] \leq_{\xi} cvr-w(\text{ver3})[\text{ver5}] \leq_{\xi} cvr-w(\text{ver3})[\text{ver4}] \leq_{\xi} cvr-w(\text{ver4})[\text{ver6}]$

- Define Versioned Objects
- Define CoVerability: a new property for distributed, versioned, objects
- Define Atomic Coverable R/W Registers (ACR)
- Compare ACR with Ranked Registers
- Use ACR to implement
 - Weak RMW objects
 - Distributed File objects

Algorithm ABD: Recalling the past

- [Attiya, Bar-Noy, Dolev 1996] (Dijkstra Prize 2011)
- Order Operations by using $\langle ts, v \rangle$ pairs.

write(v)



Servers



Writer Protocol

- $ts++$ //increment ts
- Send $\langle ts, v \rangle$ to all
- Wait for a majority to reply

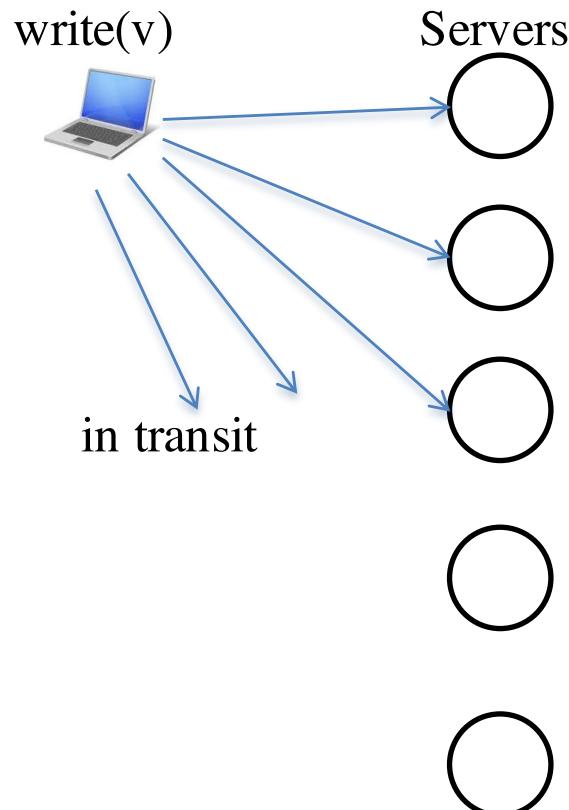
Server Protocol (Upon rcv a

$msg\ m\ from\ i$)

- if $m.ts > ts$
 - Update local $\langle ts, v \rangle$
- Send local $\langle ts, v \rangle$ to i

Algorithm ABD: Recalling the past

- [Attiya, Bar-Noy, Dolev 1996] (Dijkstra Prize 2011)
- Order Operations by using $\langle ts, v \rangle$ pairs.

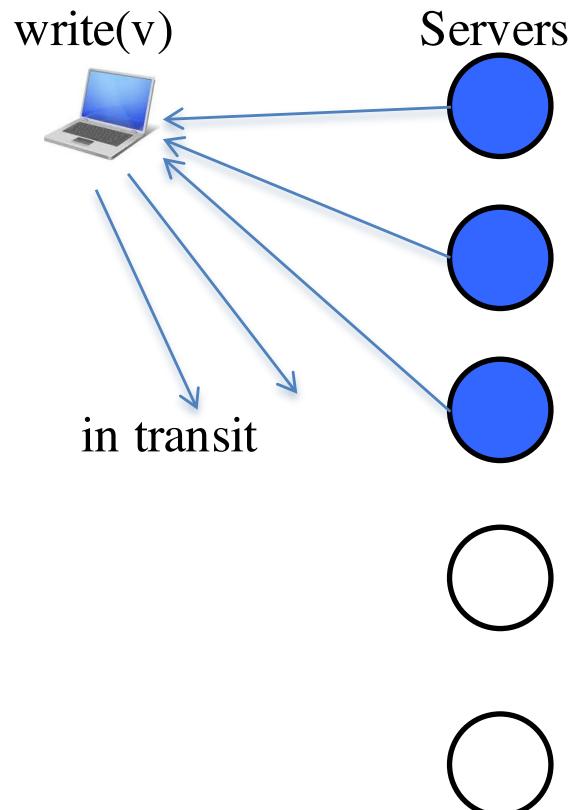


Writer Protocol

- $ts++$ //increment ts
- Send $\langle ts, v \rangle$ to all
- Wait for a majority to reply

Algorithm ABD: Recalling the past

- [Attiya, Bar-Noy, Dolev 1996] (Dijkstra Prize 2011)
- Order Operations by using $\langle ts, v \rangle$ pairs.

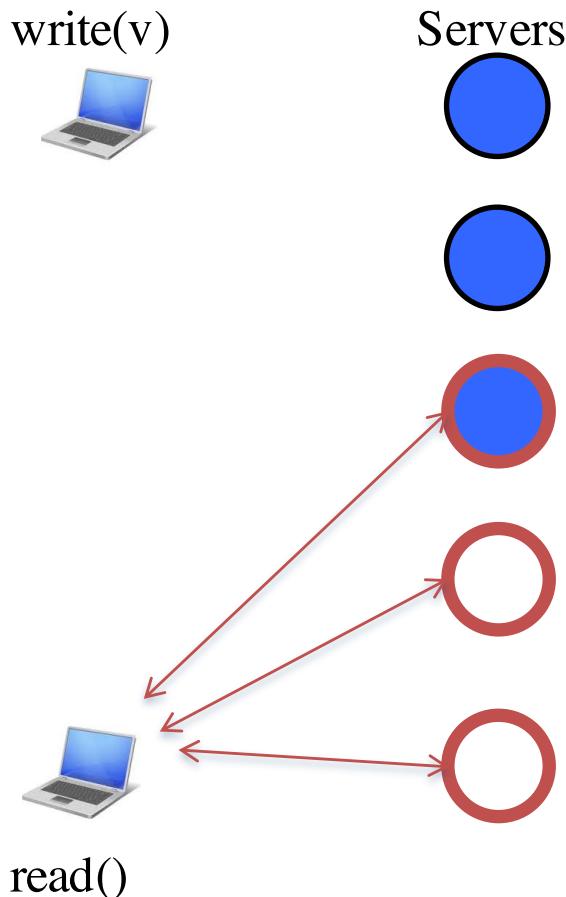


Writer Protocol

- $ts++$ //increment ts
- Send $\langle ts, v \rangle$ to all
- Wait for a majority to reply

Algorithm ABD: Recalling the past

- [Attiya, Bar-Noy, Dolev 1996] (Dijkstra Prize 2011)
- Order Operations by using $\langle ts, v \rangle$ pairs.



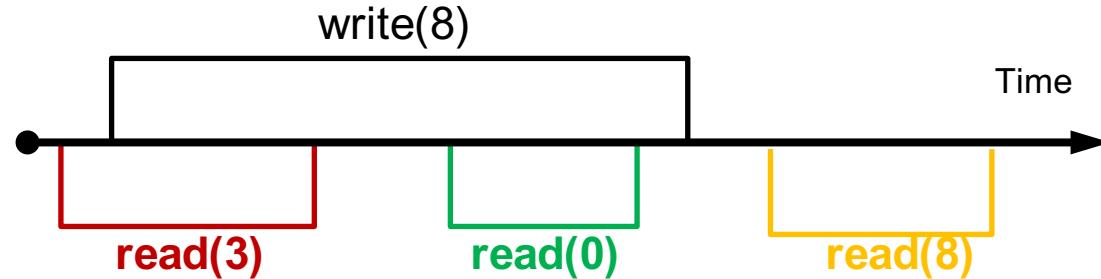
Reader Protocol (2 phases)

- Phase 1:
 - Send read to all
 - Collect $\langle ts, v \rangle$ from a majority
 - Discover $\max(\langle ts, v \rangle)$
- Phase 2:
 - Send $\max(\langle ts, v \rangle)$ to all
 - Collect ack from a majority and return v

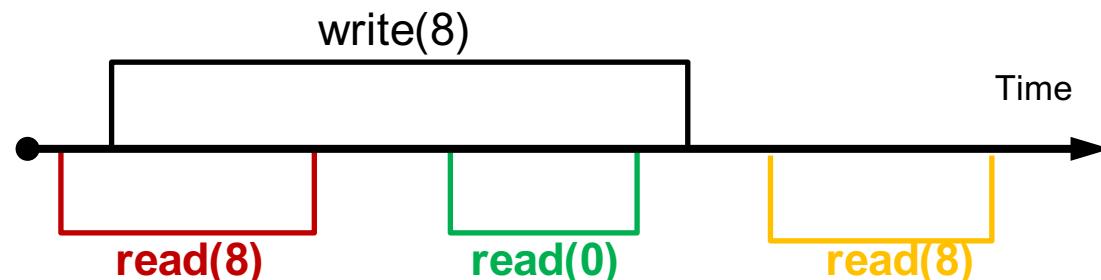
Reads must Write!

Consistency Semantics [Lamport86]

Safety



Regularity



Atomicity

