

Dependable Distributed Systems

Master of Science in Engineering in Computer Science

AA 2021/2022

LECTURE 21: REGISTERS IN PRESENCE OF BYZANTINE
PROCESSES

Safe Register Specification

Module 4.5: Interface and properties of a $(1, N)$ Byzantine safe register

Module:

Name: $(1, N)$ -ByzantineSafeRegister, **instance** *bonsr*, with writer *w*.

Events:

Request: $\langle \textit{bonsr}, \textit{Read} \rangle$: Invokes a read operation on the register.

Request: $\langle \textit{bonsr}, \textit{Write} \mid v \rangle$: Invokes a write operation with value *v* on the register.
Executed only by process *w*.

Indication: $\langle \textit{bonsr}, \textit{ReadReturn} \mid v \rangle$: Completes a read operation on the register
with return value *v*.

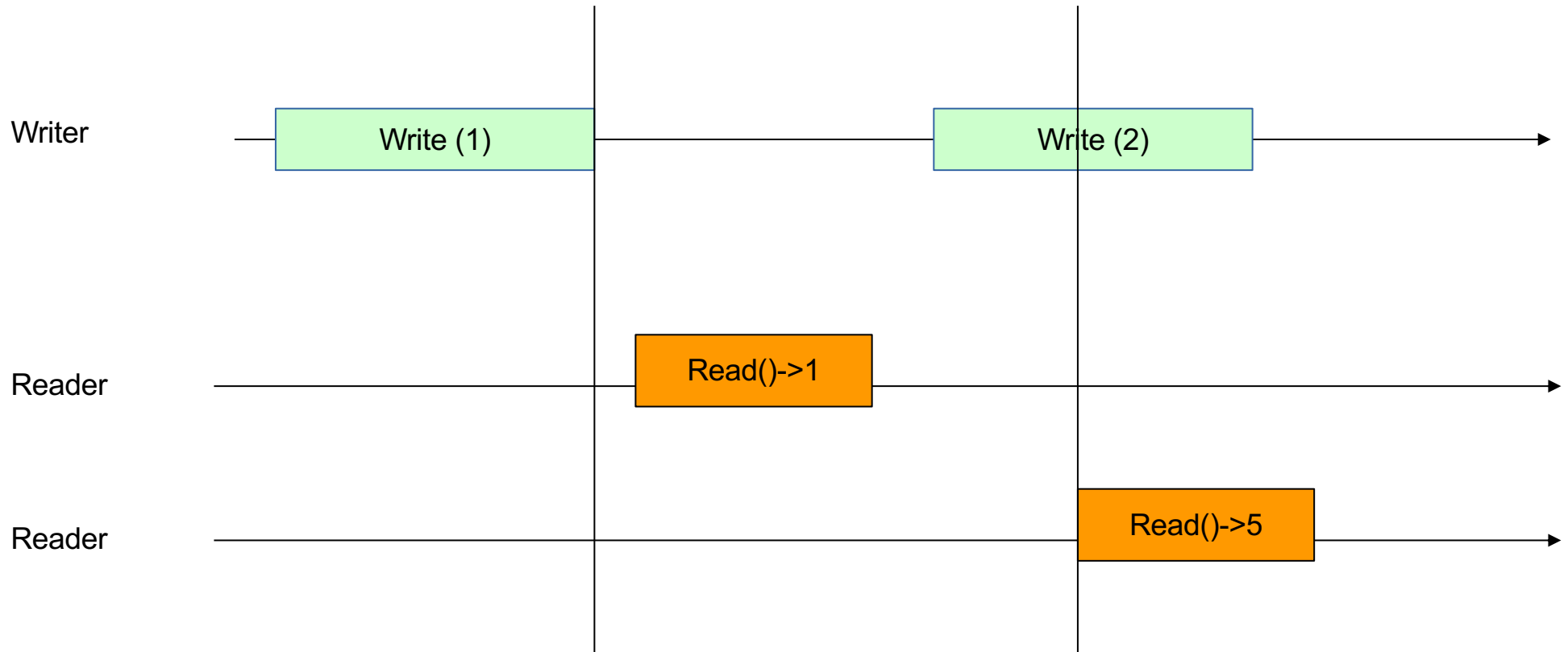
Indication: $\langle \textit{bonsr}, \textit{WriteReturn} \rangle$: Completes a write operation on the register.
Occurs only at process *w*.

Properties:

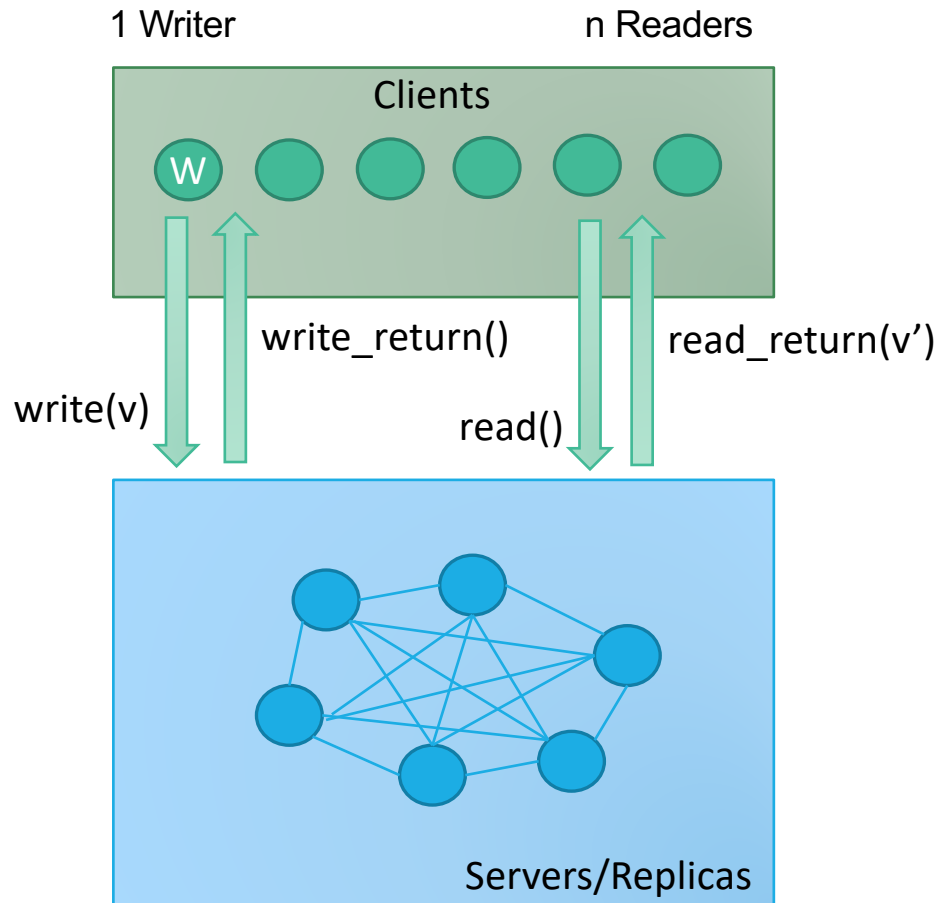
BONSR1: Termination: If a correct process invokes an operation, then the operation eventually completes.

BONSR2: Validity: A read that is not concurrent with a write returns the last value written.

Safe Register



Byzantine Tolerant Safe Register (1,n)



- Client-Server paradigm
- Asynchronous System
- authenticated perfect point-to-point link
- f servers may be Byzantine
- any client may crash but they cannot be Byzantine

Do not confuse
 n with N

Safe Register Intuition

We have to assure that once writer returns from a write operation, then any following read operation returns the last written value.

- **Write operation:** sends $\langle v, wts \rangle$ to servers and waits for ACK messages.

How many ACK messages?

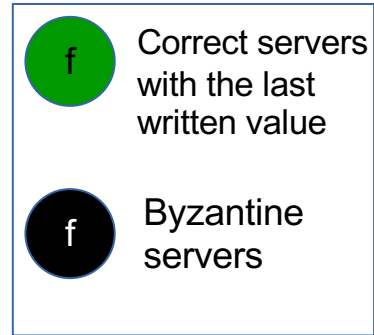
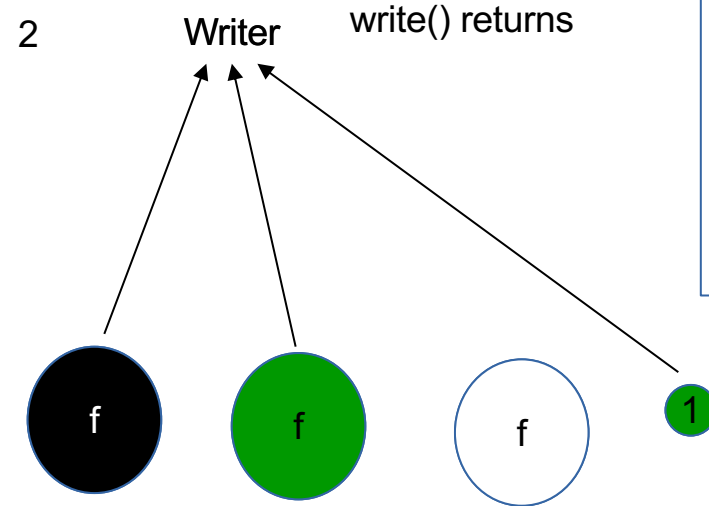
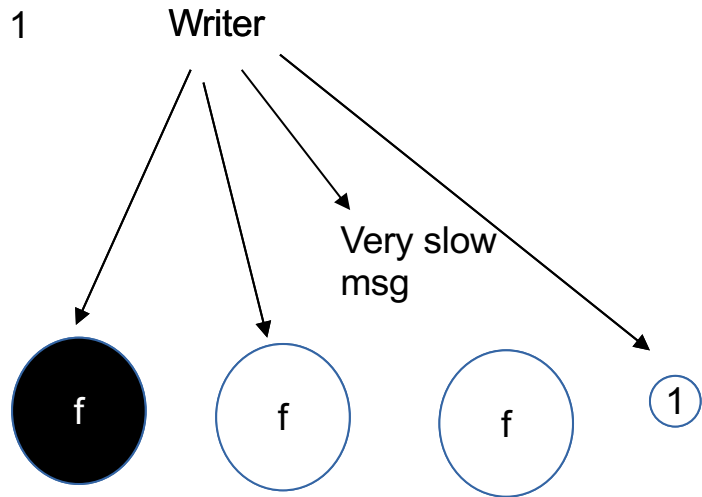
Enough to be sure that enough correct servers deliver $\langle v, wts \rangle$

- **Read operation:** sends a read request and waits for reply messages.

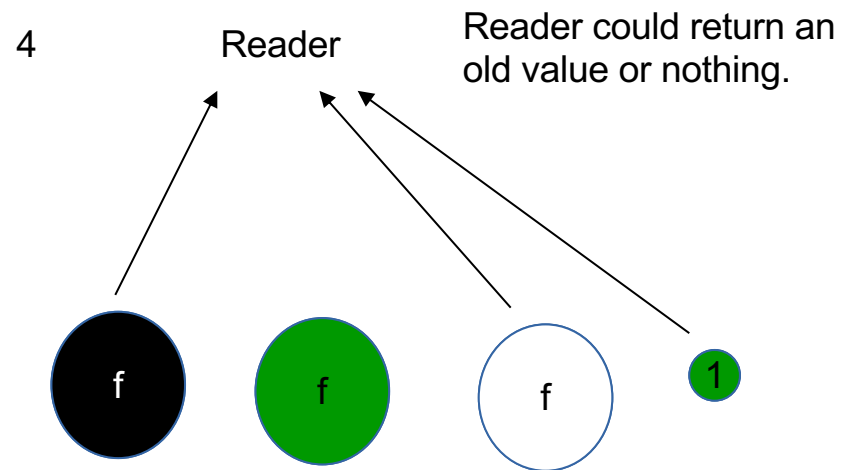
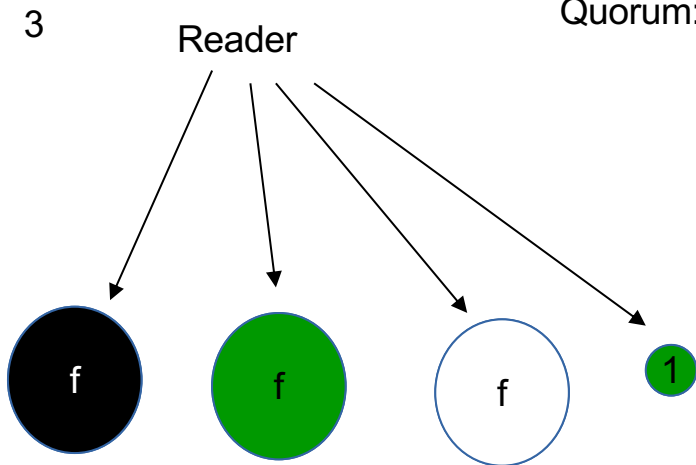
How many reply messages?

Enough to be able to read newest value, not an old value or never written

How large should the quorum be?



If we consider: $N > 3f$
Quorum: $(N+f)/2$



Masking Quorum

The kind of quorum working for Byzantine Broadcast here is not enough.

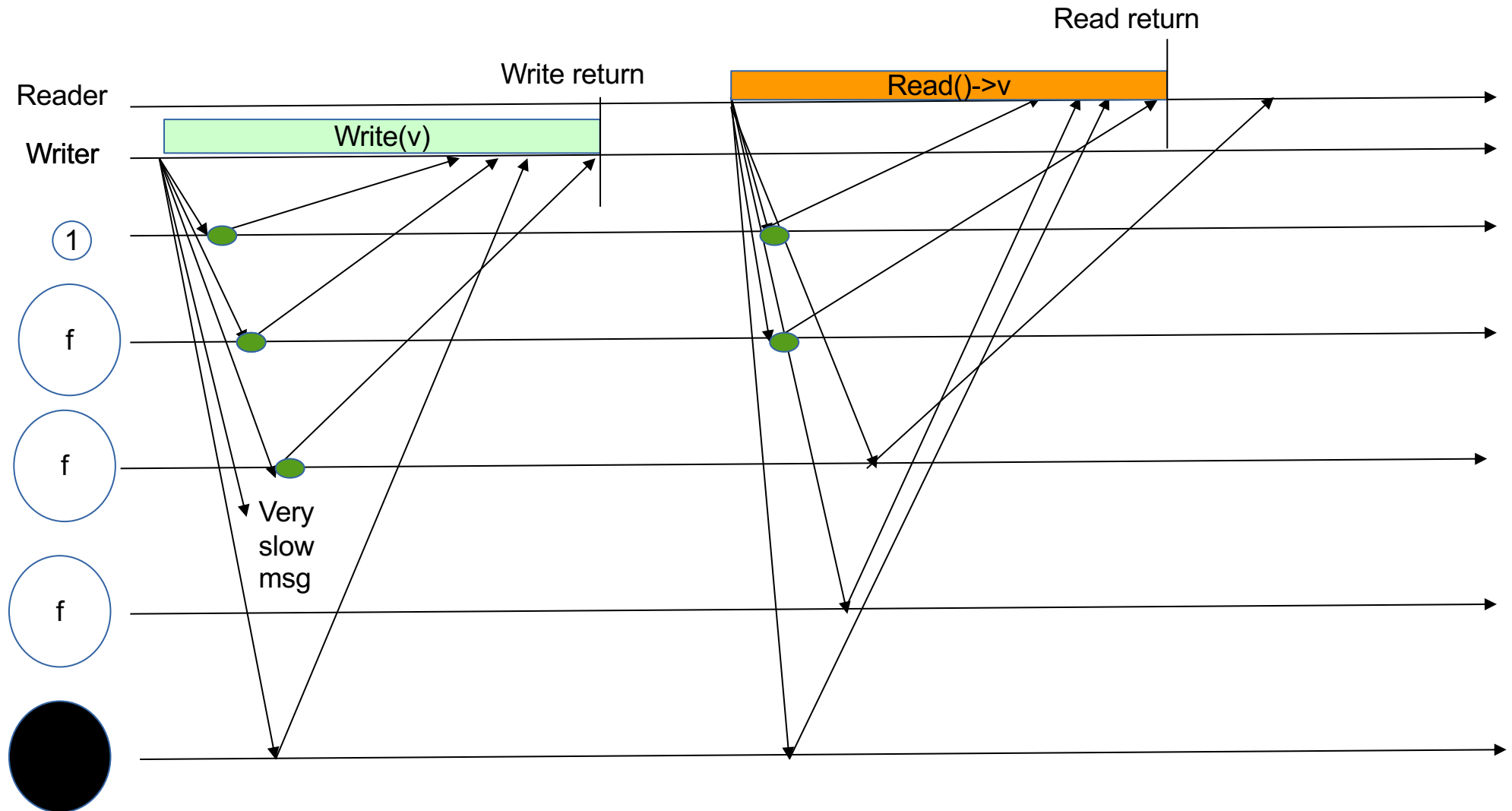
Safe Register has a stronger semantic with respect to broadcast. We require that a write operation is visible to all once it terminates.

To implement safe registers we use **Masking Quorums**:

$$N > 4f$$

$$\text{Quorum: } (N+2f)/2 \text{ (i.e., } 3f+1)$$

$N > 4f$
Quorum: $(N+2f)/2$



Algorithm 4.14: Byzantine Masking Quorum

Implements:

$(1, N)$ -ByzantineSafeRegister, **instance** *bonsr*, with writer *w*.

Uses:

AuthPerfectPointToPointLinks, **instance** *al*.

upon event $\langle \textit{bonsr}, \textit{Init} \rangle$ **do**

$(ts, val) := (0, \perp)$;

$wt_s := 0$;

$acklist := [\perp]^N$;

$rid := 0$;

$readlist := [\perp]^N$;

upon event $\langle \textit{bonsr}, \textit{Write} \mid v \rangle$ **do**

$wt_s := wt_s + 1$;

$acklist := [\perp]^N$;

forall $q \in \Pi$ **do**

trigger $\langle \textit{al}, \textit{Send} \mid q, [\textit{WRITE}, wt_s, v] \rangle$;

// only process *w*

upon event $\langle \textit{al}, \textit{Deliver} \mid p, [\textit{WRITE}, ts', v'] \rangle$ **such that** $p = w$ **do**

if $ts' > ts$ **then**

$(ts, val) := (ts', v')$;

trigger $\langle \textit{al}, \textit{Send} \mid p, [\textit{ACK}, ts'] \rangle$;

Assumption

$N > 4f$

upon event $\langle \textit{al}, \textit{Deliver} \mid q, [\textit{ACK}, ts'] \rangle$ **such that** $ts' = wt_s$ **do**

$acklist[q] := \textit{ACK}$;

if $\#(acklist) > (N + 2f)/2$ **then**

$acklist := [\perp]^N$;

trigger $\langle \textit{bonsr}, \textit{WriteReturn} \rangle$;

upon event $\langle \textit{bonsr}, \textit{Read} \rangle$ **do**

$rid := rid + 1$;

$readlist := [\perp]^N$;

forall $q \in \Pi$ **do**

trigger $\langle \textit{al}, \textit{Send} \mid q, [\textit{READ}, rid] \rangle$;

upon event $\langle \textit{al}, \textit{Deliver} \mid p, [\textit{READ}, r] \rangle$ **do**

trigger $\langle \textit{al}, \textit{Send} \mid p, [\textit{VALUE}, r, ts, val] \rangle$;

upon event $\langle \textit{al}, \textit{Deliver} \mid q, [\textit{VALUE}, r, ts', v'] \rangle$ **such that** $r = rid$ **do**

$readlist[q] := (ts', v')$;

if $\#(readlist) > \frac{N+2f}{2}$ **then**

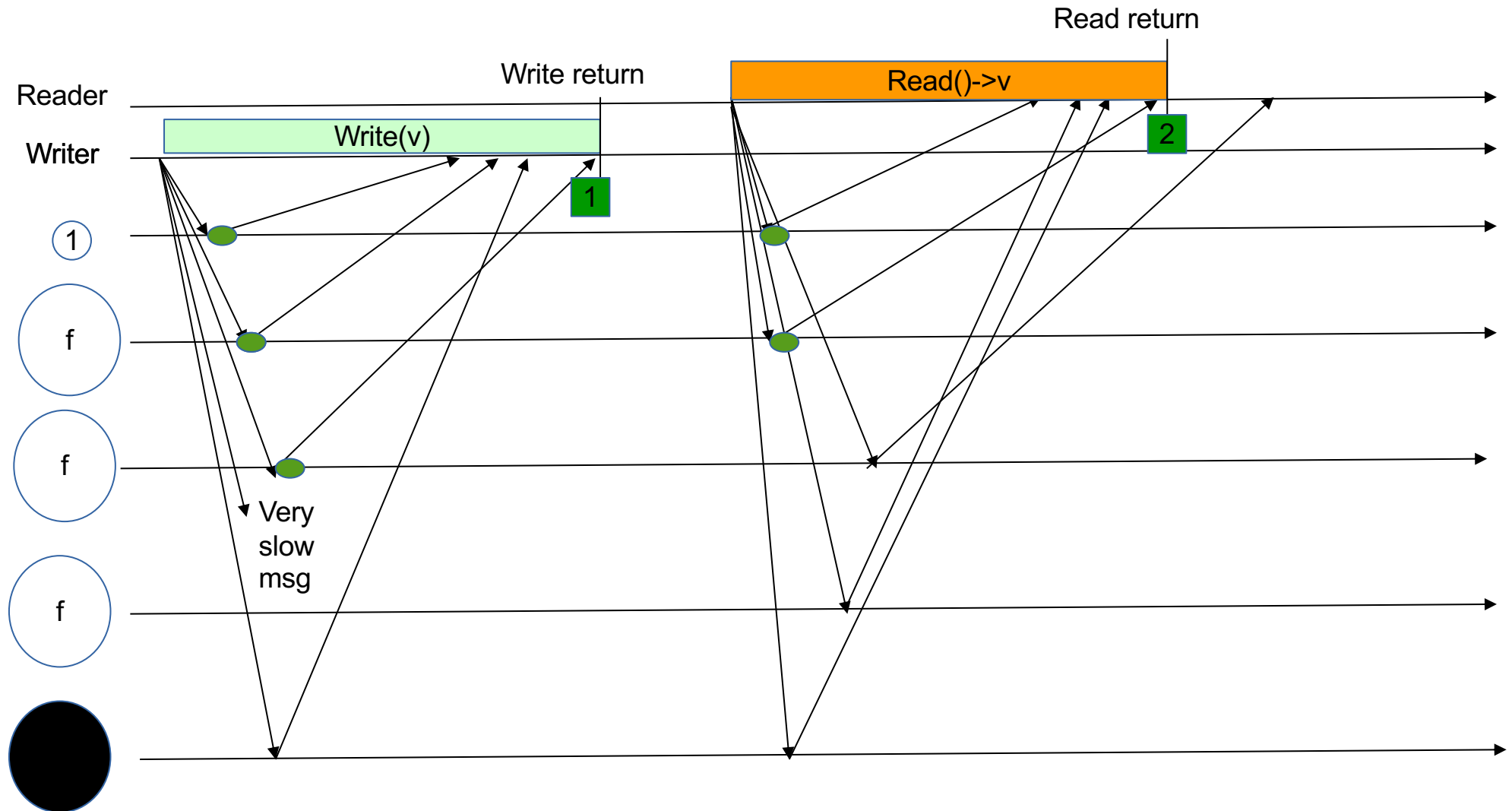
$v := \textit{byzhighestval}(readlist)$;

$readlist := [\perp]^N$;

trigger $\langle \textit{bonsr}, \textit{ReadReturn} \mid v \rangle$;

byzhighestval (\cdot): selects the value from the pair that occurs more than f time and with the highest timestamp. If no pair exists, the reader selects a default value v_0 from the domain of the register.

$N > 4f$
Quorum: $(N+2f)/2$



Regular Registers

The specification does not change

We will discuss two implementations:

1. Using cryptography
2. Without cryptography

Regular Register

Implementation with cryptographic assumptions

BASIC IDEA -> evolution of Majority voting Algorithm

- the writer signs the timestamp/value pair
- Processes store it together with the signature
- The reader verifies the signature on each timestamp/value pair received in a VALUE message and ignores those with invalid signatures



A Byzantine process is prevented from returning an arbitrary timestamp value in the VALUE message, although it may include a signed value with an outdated timestamp

Regular Register

Implementation with cryptographic assumptions

Algorithm 4.15: Authenticated-Data Byzantine Quorum

Implements:

$(1, N)$ -ByzantineRegularRegister, **instance** *bonrr*, with writer *w*.

Uses:

AuthPerfectPointToPointLinks, **instance** *al*.

upon event $\langle \text{bonrr}, \text{Init} \rangle$ **do**

$(ts, val, \sigma) := (0, \perp, \perp);$

$wts := 0;$

$acklist := [\perp]^N;$

$rid := 0;$

$readlist := [\perp]^N;$

upon event $\langle \text{bonrr}, \text{Write} \mid v \rangle$ **do**

$wts := wts + 1;$

$acklist := [\perp]^N;$

$\sigma := \text{sign}(\text{self}, \text{bonrr} \parallel \text{self} \parallel \text{WRITE} \parallel wts \parallel v);$

forall $q \in \Pi$ **do**

trigger $\langle al, \text{Send} \mid q, [\text{WRITE}, wts, \sigma] \rangle;$

// only process *w*

upon event $\langle al, \text{Deliver} \mid p, [\text{WRITE}, ts', v', \sigma'] \rangle$ **such that** $p = w$ **do**

if $ts' > ts$ **then**

$(ts, val, \sigma) := (ts', v', \sigma');$

trigger $\langle al, \text{Send} \mid p, [\text{ACK}, ts'] \rangle;$

upon event $\langle al, \text{Deliver} \mid q, [\text{ACK}, ts'] \rangle$ **such that** $ts' = wts$ **do**

$acklist[q] := \text{ACK};$

if $\#(acklist) < (N + f)/2$ **then**

$acklist := [\perp]^N;$

trigger $\langle \text{bonrr}, \text{WriteReturn} \rangle;$

i.e. $> 2f$

Assumption
 $N > 3f$

upon event $\langle \text{bonrr}, \text{Read} \rangle$ **do**

$rid := rid + 1;$

$readlist := [\perp]^N;$

forall $q \in \Pi$ **do**

trigger $\langle al, \text{Send} \mid q, [\text{READ}, rid] \rangle;$

upon event $\langle al, \text{Deliver} \mid p, [\text{READ}, r] \rangle$ **do**

trigger $\langle al, \text{Send} \mid p, [\text{VALUE}, r, ts, val, \sigma] \rangle;$

upon event $\langle al, \text{Deliver} \mid q, [\text{VALUE}, r, ts', v', \sigma'] \rangle$ **such that** $r = rid$ **do**

if $\text{verifysig}(q, \text{bonrr} \parallel w \parallel \text{WRITE} \parallel ts' \parallel v', \sigma')$ **then**

$readlist[q] := (ts', v');$

if $\#(readlist) > \frac{N+f}{2}$ **then**

$v := \text{highestval}(readlist);$

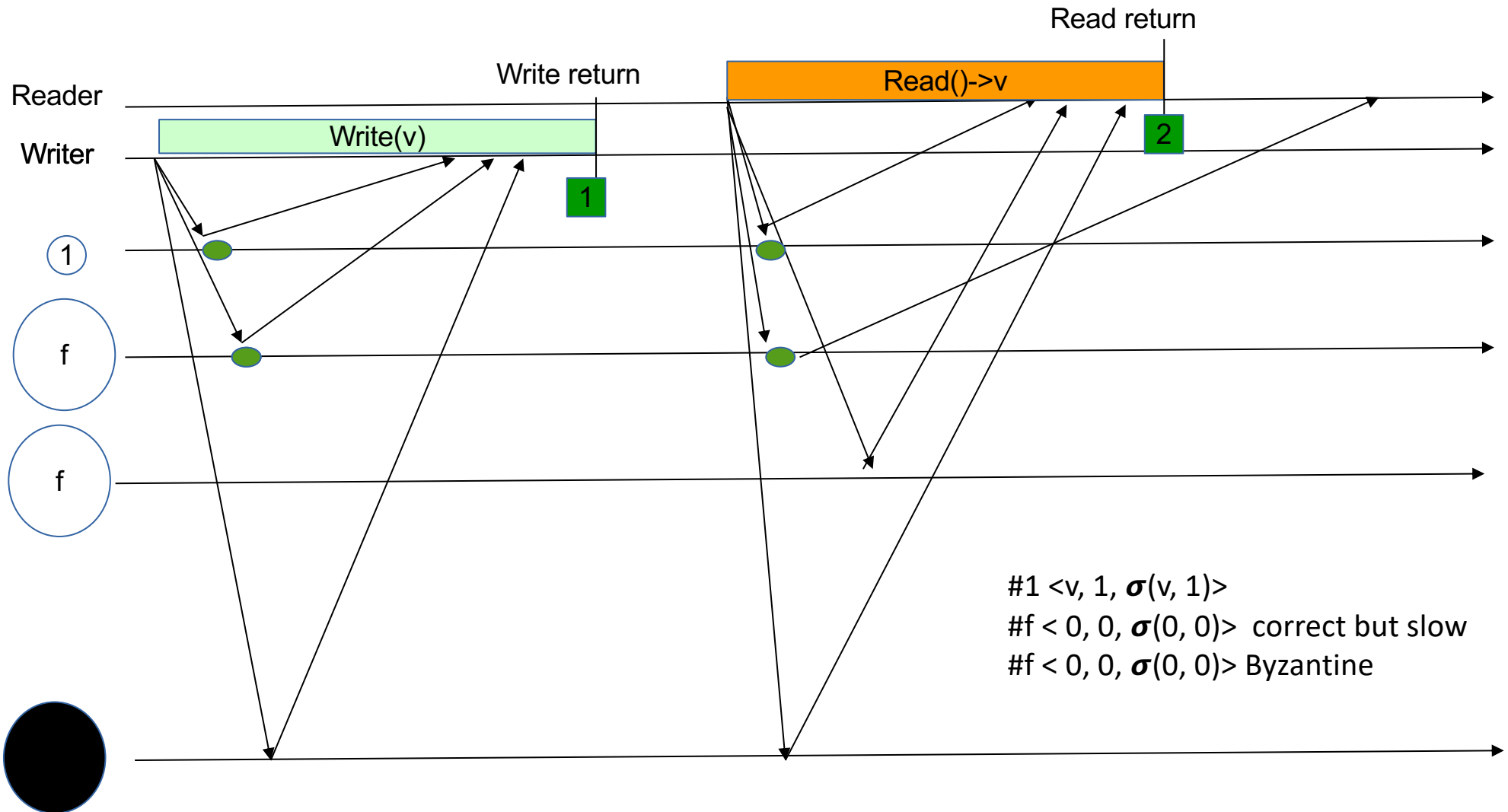
$readlist := [\perp]^N;$

trigger $\langle \text{bonrr}, \text{ReadReturn} \mid v \rangle;$

i.e. $> 2f$

$N > 3f$
 Quorum: $(N+2f)/2$

Very
slow
msg

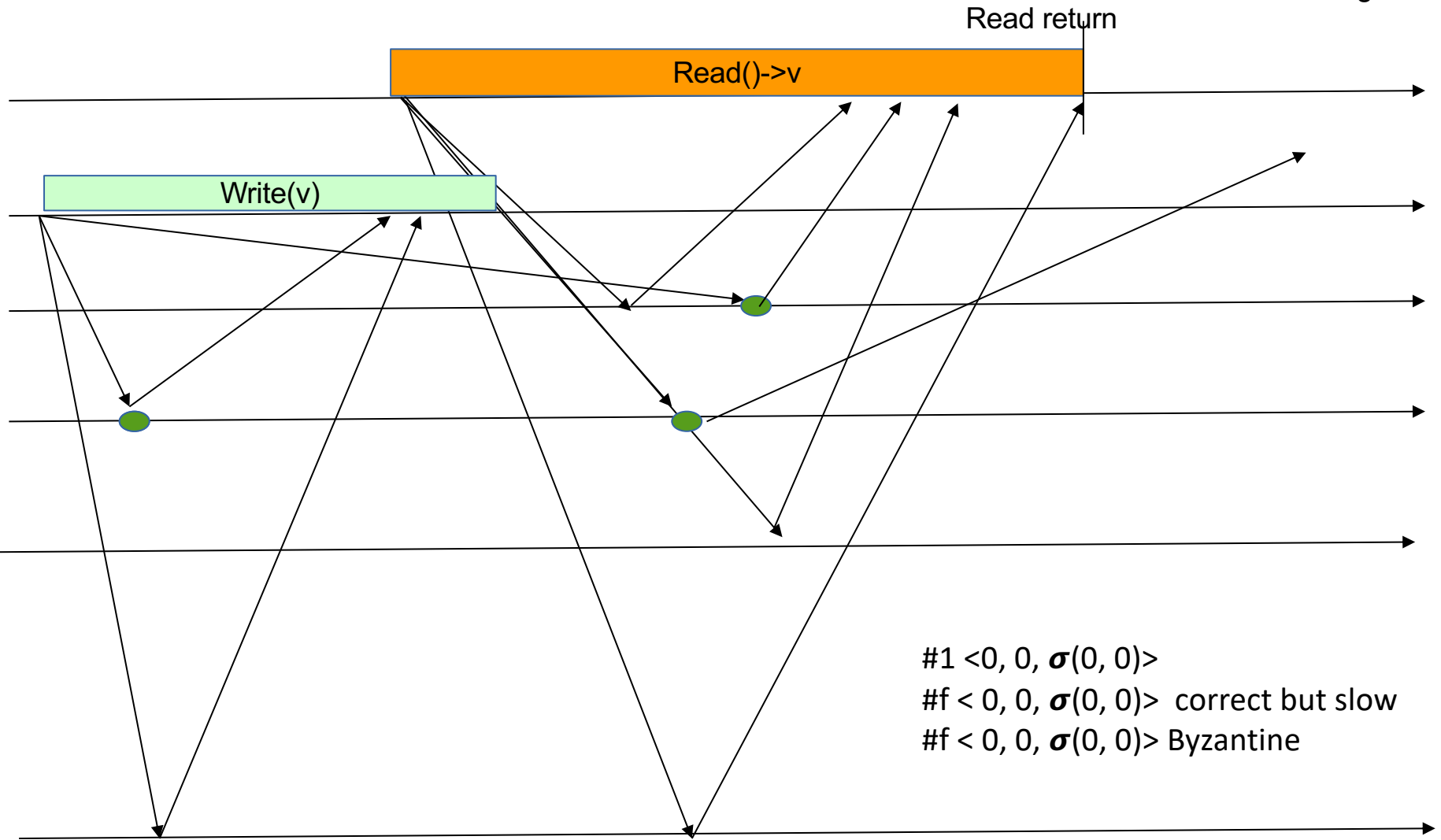
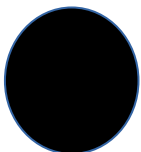


$N > 3f$
Quorum: $(N+2f)/2$

Reader

Writer

1
f
f



Regular Register Implementation without cryptographic assumptions

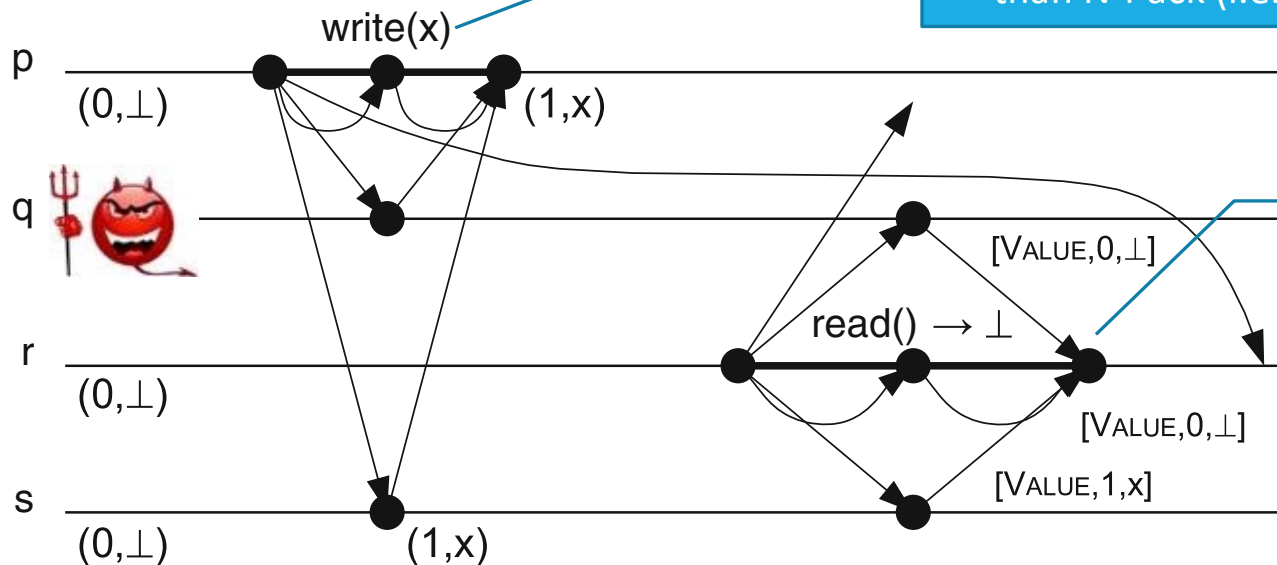
the writer process p uses two phases to write a new

- a pre-write phase and
- a write phase

Assumption
 $N > 3f$

Why one phase is not enough?

In order to terminate, the write cannot wait more than $N - f$ ack (i.e., 3)



The reader is not able to choose a value:

1. Value \perp could be old given that process s provided value x with ts 1
2. Value x could be a fake value generated by a Byzantine process

Regular Register

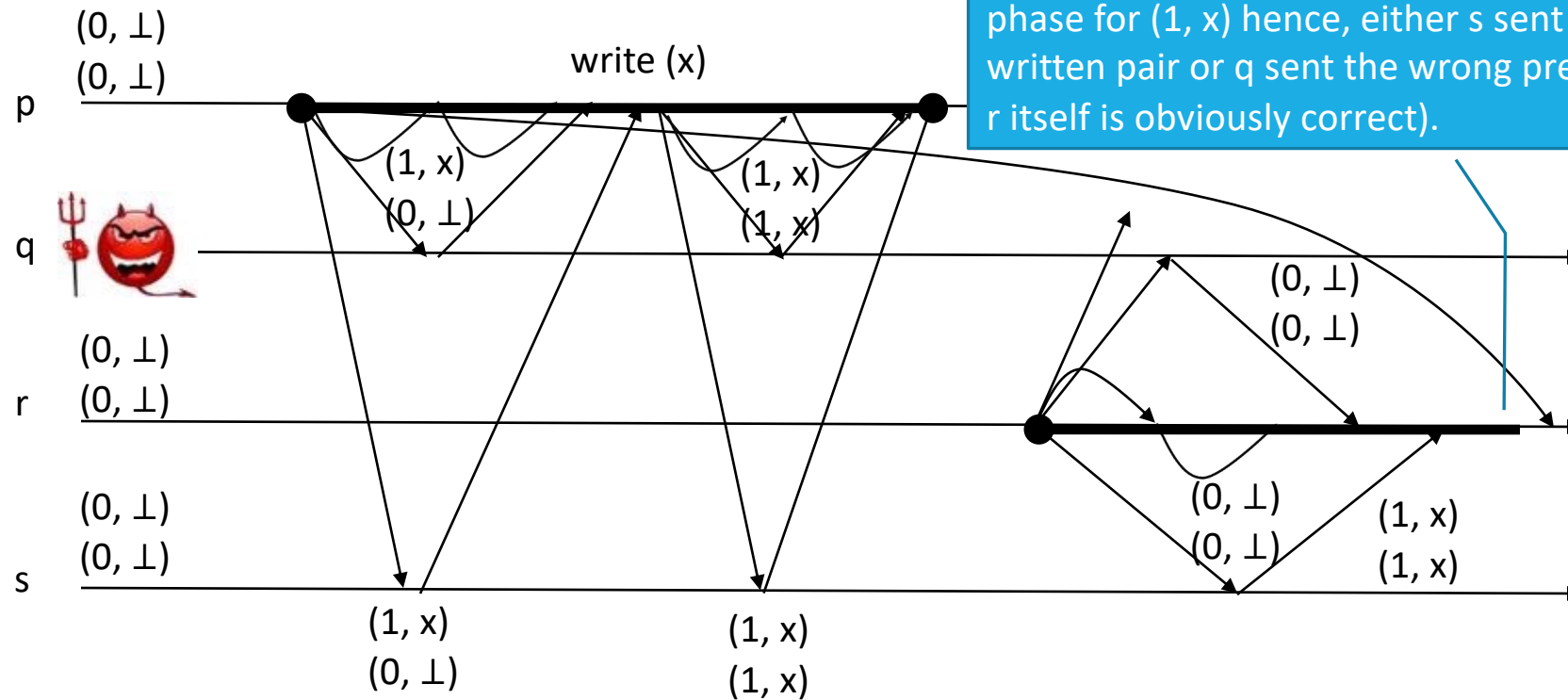
Implementation without cryptographic assumptions

the writer process p uses two phases to write a new

- a **pre-write phase**: the writer sends PREWRITE messages with the current timestamp/value pair. Then it waits until it receives PREACK messages from $N - f$ processes
- a **write phase**: the writer sends ordinary WRITE messages, again containing the current timestamp/value pair, and then waits until it receives ACK messages from $N - f$ processes

Every process stores two timestamp/value pairs, one from the pre-write phase and one from the write phase

Intuition



Regular Register

write Implementation without cryptographic assumptions

Algorithm 4.16: Double-Write Byzantine Quorum (part 1, write)

Implements:

$(1, N)$ -ByzantineRegularRegister, **instance** *bonrr*, with writer *w*.

Uses:

AuthPerfectPointToPointLinks, **instance** *al*.

upon event $\langle \text{bonrr}, \text{Init} \rangle$ **do**

$(pts, pval) := (0, \perp);$

$(ts, val) := (0, \perp);$

$(wts, wval) := (0, \perp);$

$preacklist := [\perp]^N;$

$acklist := [\perp]^N;$

$rid := 0;$

$readlist := [\perp]^N;$

upon event $\langle \text{bonrr}, \text{Write} \mid v \rangle$ **do**

$(wts, wval) := (wts + 1, v);$

$preacklist := [\perp]^N;$

$acklist := [\perp]^N;$

forall $q \in \Pi$ **do**

trigger $\langle al, \text{Send} \mid q, [\text{PREWRITE}, wts, wval] \rangle;$

// only process *w*

upon event $\langle al, \text{Deliver} \mid p, [\text{PREWRITE}, pts', pval'] \rangle$

such that $p = w \wedge pts' = pts + 1$ **do**

$(pts, pval) := (pts', pval');$

trigger $\langle al, \text{Send} \mid p, [\text{PREACK}, pts] \rangle;$

upon event $\langle al, \text{Deliver} \mid q, [\text{PREACK}, pts'] \rangle$ **such that** $pts' = wts$ **do**

$preacklist[q] := \text{PREACK};$

if $\#(preacklist) \geq N - f$ **then**

$preacklist := [\perp]^N;$

forall $q \in \Pi$ **do**

trigger $\langle al, \text{Send} \mid q, [\text{WRITE}, wts, wval] \rangle;$

upon event $\langle al, \text{Deliver} \mid p, [\text{WRITE}, ts', val'] \rangle$

such that $p = w \wedge ts' = pts \wedge ts' > ts$ **do**

$(ts, val) := (ts', val');$

trigger $\langle al, \text{Send} \mid p, [\text{ACK}, ts] \rangle;$

upon event $\langle al, \text{Deliver} \mid q, [\text{ACK}, ts'] \rangle$ **such that** $ts' = wts$ **do**

$acklist[q] := \text{ACK};$

if $\#(acklist) \geq N - f$ **then**

$acklist := [\perp]^N;$

trigger $\langle \text{bonrr}, \text{WriteReturn} \rangle;$

Regular Register

read Implementation without cryptographic assumptions

Algorithm 4.17: Double-Write Byzantine Quorum (part 2, read)

upon event $\langle \text{bonrr}, \text{Read} \rangle$ **do**
 $\text{rid} := \text{rid} + 1$;
 $\text{readlist} := [\perp]^N$;
 forall $q \in \Pi$ **do**
 trigger $\langle \text{al}, \text{Send} \mid q, [\text{READ}, \text{rid}] \rangle$;

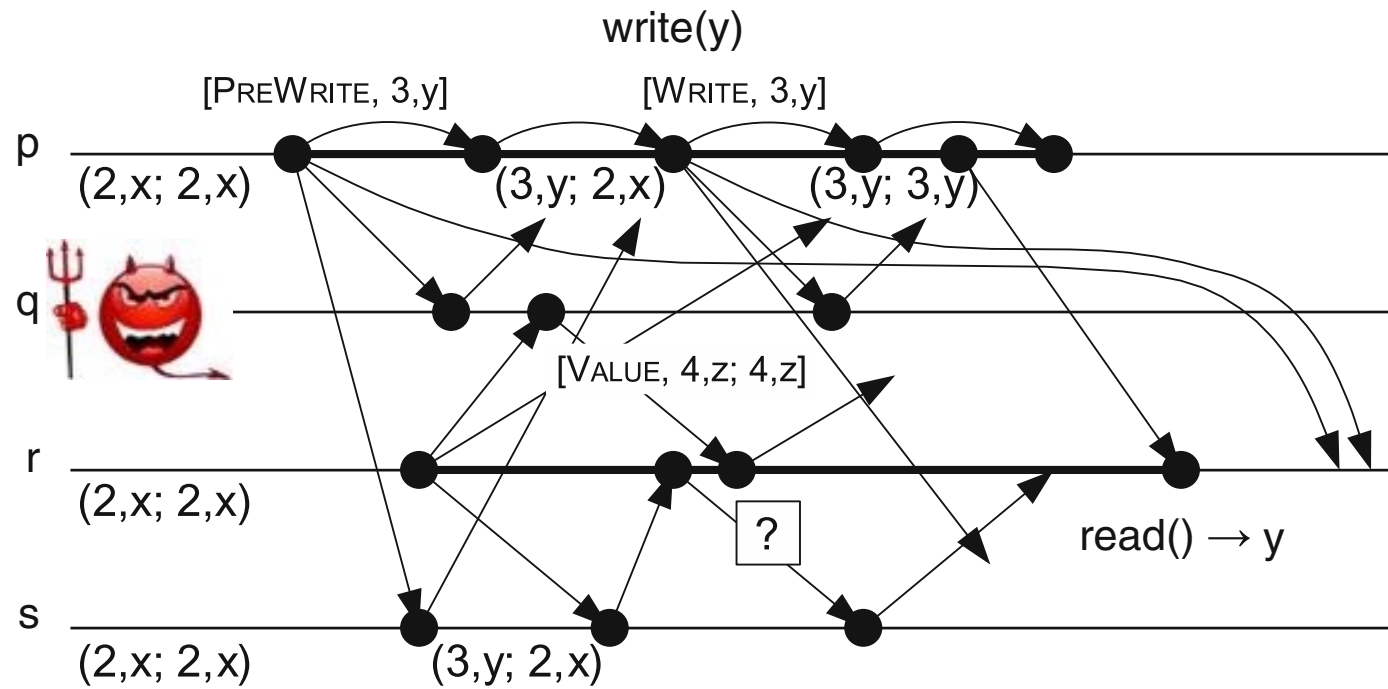
upon event $\langle \text{al}, \text{Deliver} \mid p, [\text{READ}, r] \rangle$ **do**
 trigger $\langle \text{al}, \text{Send} \mid p, [\text{VALUE}, r, \text{pts}, \text{pval}, \text{ts}, \text{val}] \rangle$;

upon event $\langle \text{al}, \text{Deliver} \mid q, [\text{VALUE}, r, \text{pts}', \text{pval}', \text{ts}', \text{val}'] \rangle$ **such that** $r = \text{rid}$ **do**
 if $\text{pts}' = \text{ts}' + 1 \vee (\text{pts}', \text{pval}') = (\text{ts}', \text{val}')$ **then**
 $\text{readlist}[q] := (\text{pts}', \text{pval}', \text{ts}', \text{val}')$;
 if exists (ts, v) in an entry of readlist **such that** $\text{authentic}(\text{ts}, v, \text{readlist}) = \text{TRUE}$
 and exists $Q \subseteq \text{readlist}$ **such that**
 $\#(Q) > \frac{N+f}{2} \wedge \text{selectedmax}(\text{ts}, v, Q) = \text{TRUE}$ **then**
 $\text{readlist} := [\perp]^N$;
 trigger $\langle \text{bonrr}, \text{ReadReturn} \mid v \rangle$;
 else
 trigger $\langle \text{al}, \text{Send} \mid q, [\text{READ}, r] \rangle$;

TRUE if readlist contains a pair (ts, v) that is found in the entries of more than f processes

readlist contains a pair (ts, v) such there is a Byzantine quorum (Q) of entries in readlist whose *highest* timestamp/value pair, selected among the pre-written or written pair of the entries, is (ts, v)

Example



Regular Register without cryptographic assumptions: Observations

Termination property must be relaxed in to *finite-write termination*

- Instead of requiring that *every* operation of a correct process eventually terminates, a *read* operation that is concurrent with *infinitely* many write operations may not terminate

It has been shown that such a relaxation is necessary

Exercise

Does the Regular Register implementation without cryptographic assumption satisfy also the atomic specification?

References

C. Cachin, R. Guerraoui and L. Rodrigues. Introduction to Reliable and Secure Distributed Programming, Springer, 2011

- Chapter 4 – Sections 4.6 and 4.7