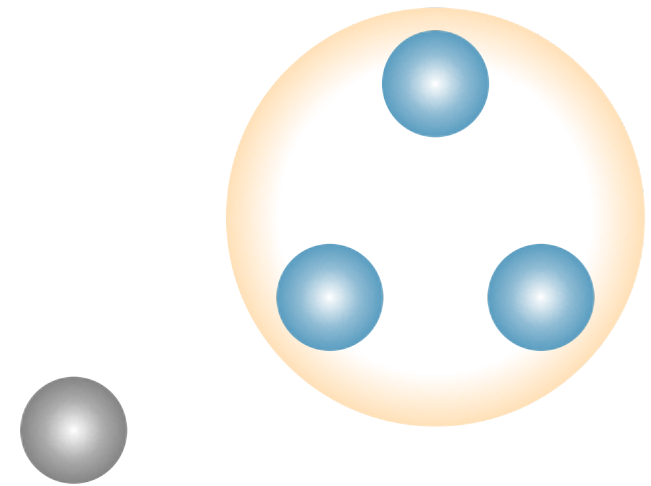# Syllabus

‣ Introduction to fault-tolerant distributed systems

‣ Models of distributed systems and related faults

‣ **Data replication**

‣ Distributed consensus

‣ State machine replication
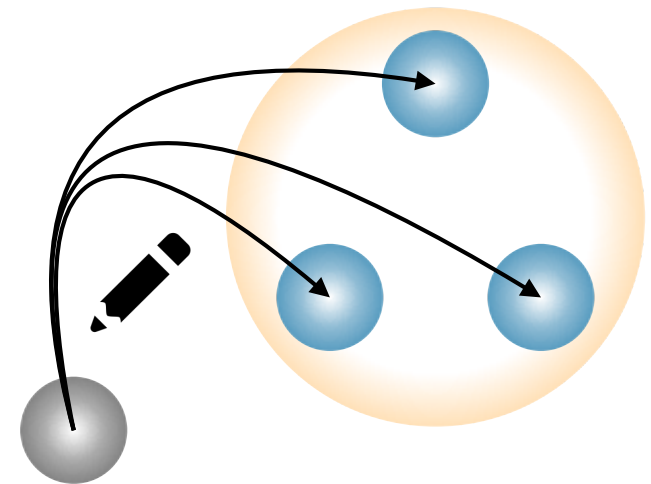
‣ Database replication

# Data replication
## ROWA

‣ Let us consider the replication of a basic storage service

  ‣ Elementary atomic read and write operations (mutable servers)

  ‣ Multiple sequential client and server processes | replicas (concurrency model)

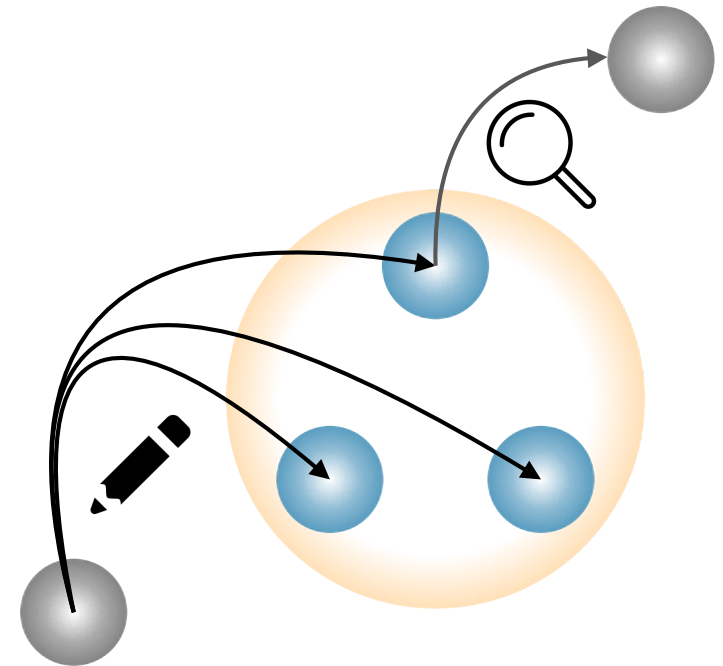  ‣ Crash faults (fault model)

# Data replication
## ROWA

‣ Let us consider the replication of a basic storage service

  ‣ Elementary atomic read and write operations (mutable servers)

  ‣ Multiple sequential client and server processes I replicas (concurrency model)

  ‣ Crash faults (fault model)

‣ Read-One-Write-All (ROWA)

‣ Let us consider the replication of a basic storage service

- ‣ Elementary atomic read and write operations (mutable servers)

- ‣ Multiple sequential client and server processes I replicas (concurrency model)

- ‣ Crash faults (fault model)

‣ Read-One-Write-All (ROWA)

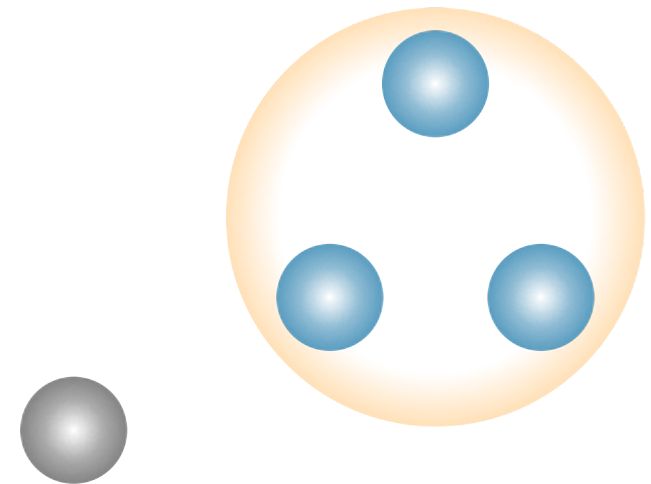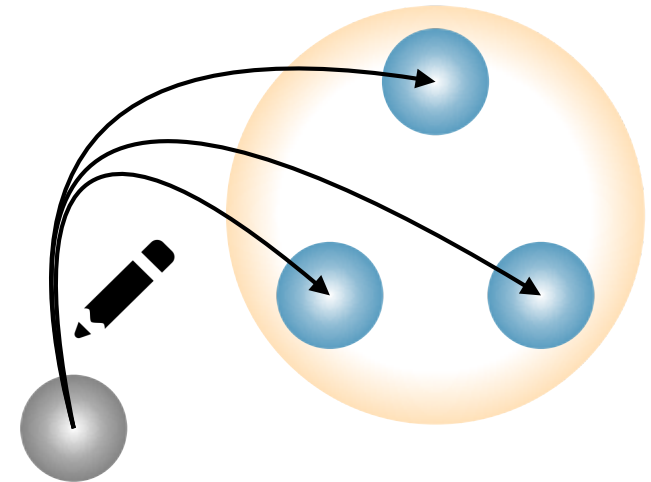- ‣ Write *value* to all replicas

# Data replication

▸ Let us consider the replication of a basic storage service

   ‣ Elementary atomic read and write operations (mutable servers)

   ‣ Multiple sequential client and server processes | replicas (concurrency model)

   ‣ Crash faults (fault model)

▸ Read-One-Write-All (ROWA)

   ‣ Write *value* to all replicas

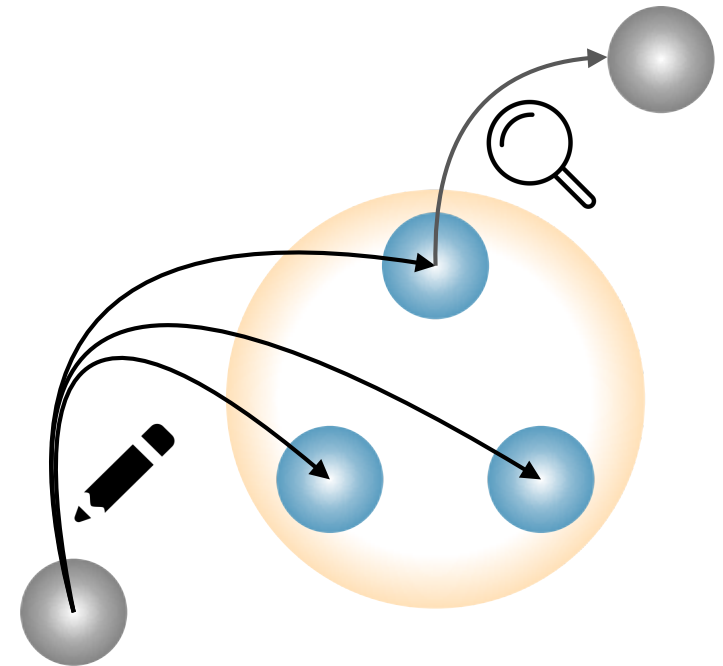   ‣ Read can be performed on any replica

# Data replication
## ROWA

‣ Let us consider the replication of a basic storage service

- ‣ Elementary atomic read and write operations (mutable servers)

- ‣ Multiple sequential client and server processes / replicas (concurrency model)

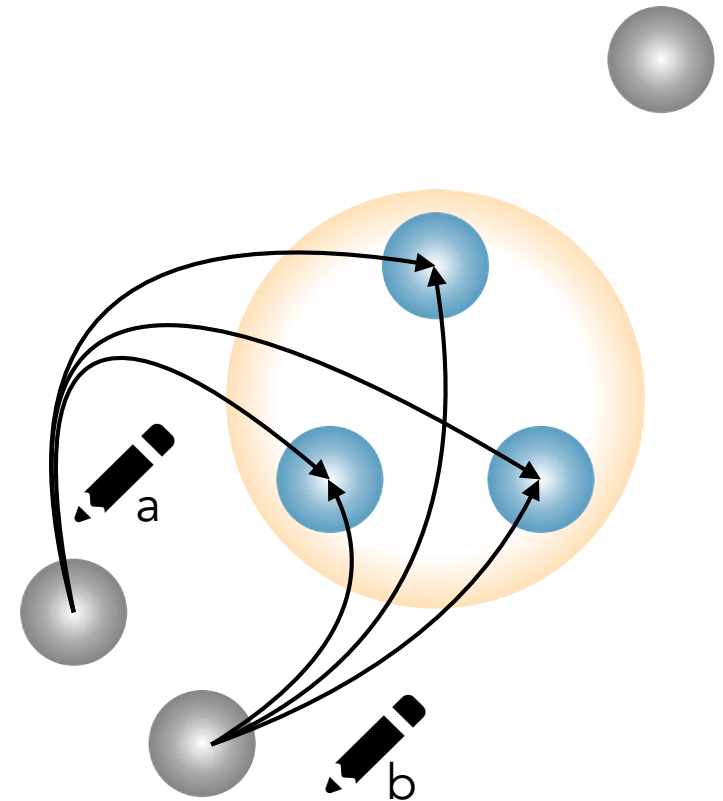- ‣ No faults (fault model)

# Data replication
## ROWA

‣ Let us consider the replication of a basic storage service

  ‣ Elementary atomic read and write operations (mutable servers)

  ‣ Multiple sequential client and server processes / replicas (concurrency model)

  ‣ No faults (fault model)
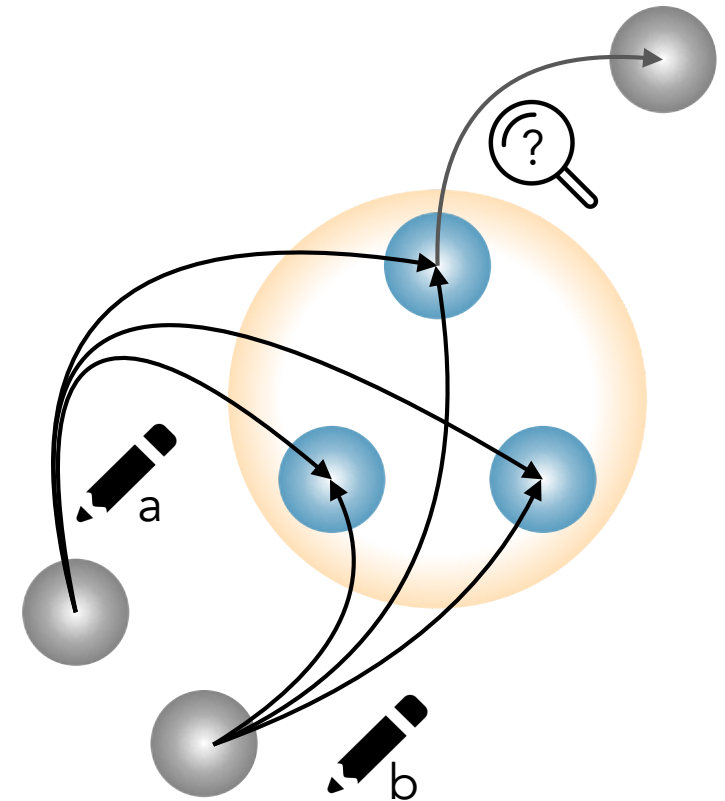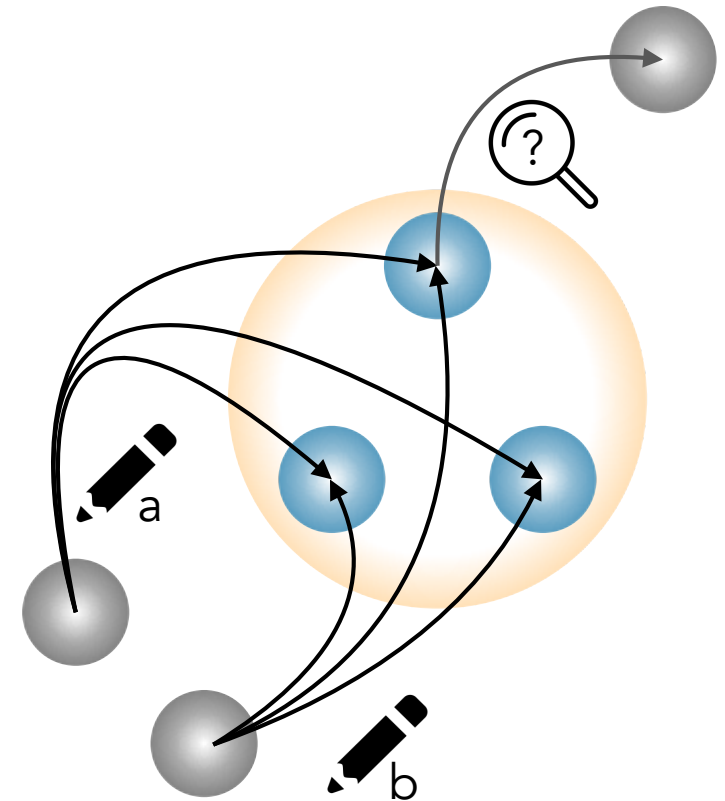
‣ Read-One-Write-All (ROWA)

# Data replication

‣ Let us consider the replication of a basic storage service

  ‣ Elementary atomic read and write operations (mutable servers)

  ‣ Multiple sequential client and server processes / replicas (concurrency model)

  ‣ No faults (fault model)

‣ Read-One-Write-All (ROWA)
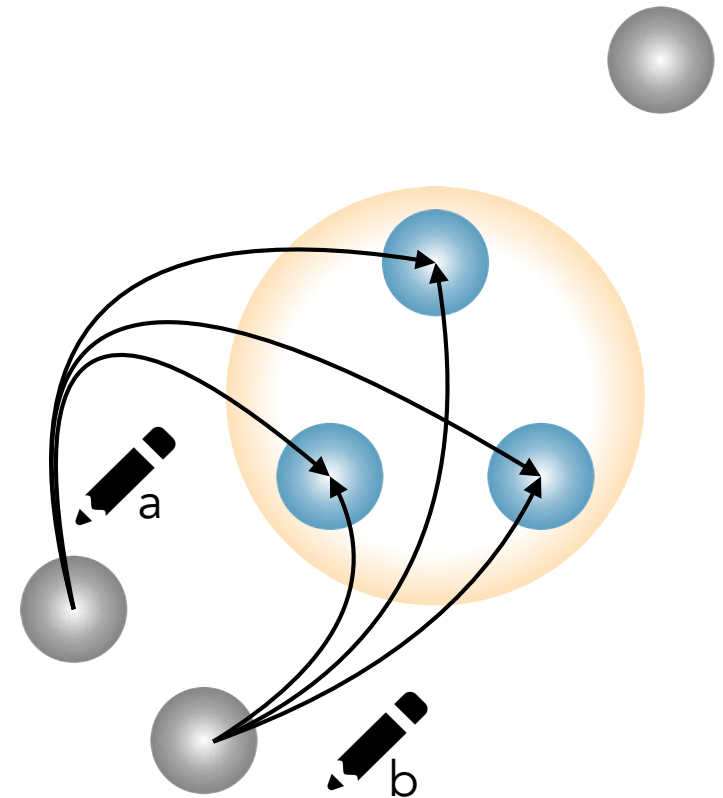
  ‣ Write *value* to all replicas

‣ Let us consider the replication of a basic storage service

- ‣ Elementary atomic read and write operations (mutable servers)

- ‣ Multiple sequential client and server processes / replicas (concurrency model)

- ‣ No faults (fault model)

‣ Read-One-Write-All (ROWA)

- ‣ Write *value* to all replicas

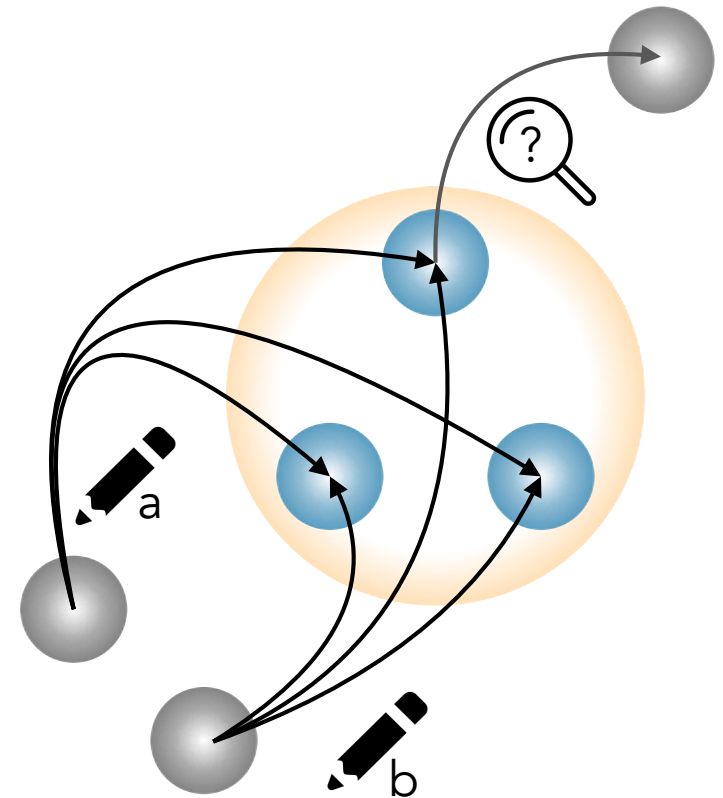- ‣ Read can be performed on any replica

# Data replication

‣ Let us consider the replication of a basic storage service

   ‣ Elementary atomic read and write operations (mutable servers)

   ‣ Multiple sequential client and server processes / replicas (concurrency model)

   ‣ No faults (fault model)

‣ Read-One-Write-All (ROWA)

   ‣ Write *value* to all replicas

   ‣ Read can be performed on any replica
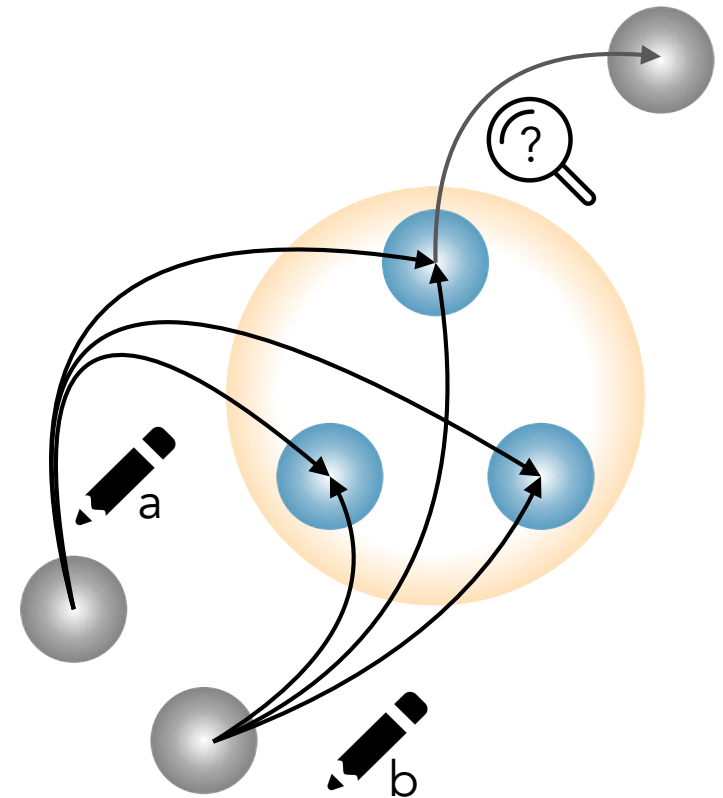
‣ Concurrent writes

# Data replication
## ROWA

‣ Let us consider the replication of a basic storage service

  ‣ Elementary atomic read and write operations (mutable servers)

  ‣ Multiple sequential client and server processes / replicas (concurrency model)

  ‣ No faults (fault model)

‣ Read-One-Write-All (ROWA)

  ‣ Write *value* to all replicas

  ‣ Read can be performed on any replica

‣ Concurrent writes

  ‣ What does a subsequent read return?

‣ Let us consider the replication of a basic storage service

  ‣ Elementary atomic read and write operations (mutable servers)

  ‣ Multiple sequential client and server processes / replicas (concurrency model)

  ‣ No faults (fault model)

‣ Read-One-Write-All (ROWA)

  ‣ Write *value* to all replicas

  ‣ Read can be performed on any replica

‣ Concurrent writes

  ‣ What does a subsequent read return?

  ‣ We may end up with replicas writing a ➔ b and others writing b ➔ a 😳

# Data replication

‣ Let us consider the replication of a basic storage service

- ‣ Elementary atomic read and write operations (mutable servers)

- ‣ Multiple sequential client and server processes / replicas (concurrency model)

- ‣ Crash faults (fault model)

‣ Read-One-Write-All (ROWA)

- ‣ Write *value* to all replicas

- ‣ Read can be performed on any replica

‣ Concurrent writes

# Data replication

‣ Let us consider the replication of a basic storage service

  ‣ Elementary atomic read and write operations (mutable servers)

  ‣ Multiple sequential client and server processes / replicas (concurrency model)

  ‣ Crash faults (fault model)

‣ Read-One-Write-All (ROWA)

  ‣ Write *value* to all replicas

  ‣ Read can be performed on any replica

‣ Concurrent writes

  ‣ What does a subsequent read return?

# Data replication
## ROWA

‣ Let us consider the replication of a basic storage service

   ‣ Elementary atomic read and write operations (mutable servers)

   ‣ Multiple sequential client and server processes / replicas (concurrency model)

   ‣ Crash faults (fault model)

‣ Read-One-Write-All (ROWA)

   ‣ Write *value* to all replicas

   ‣ Read can be performed on any replica

‣ Concurrent writes

   ‣ What does a subsequent read return?

   ‣ We may end up with replicas writing a and then b and others writing b and then a

‣ How can we prevent stale writes?

‣ Let us add a *version* to the state of the replicas

   ‣ $v_i$ is replica i's version, initially null

‣ Write operations now become more complex

At client proxies

```
write (value)
    _read (_, vmax) from some replica
    v_write (value, vmax+1) to all replicas
```

At replica i

```
_write (value, version)
    if version > vi then
        xi = value
        vi = version
```

‣ How to define version? Is a simple scalar enough?

‣ Eg. version = (counter, pid)

```
vi > vj ::
    vi.counter > vj.counter OR
    vi.counter = vj.counter AND vi.pid > vj.pid
```

‣ <u>ROWA is not fault tolerant</u>

‣ Let us introduce Quorums

    ‣ A quorum is a set of replicas

    ‣ We will refer to two quorums: a write quorum Qw and a read quorum Qr

    ‣ We will write to some Qw set of replicas and will read from some Qr set

    ‣ ROWA is a particular case in which $|Qw| = n$ and $|Qr| = 1$

‣ If we make $|Qw| < n$ then the replicated system becomes <u>fault tolerant</u>

‣ Quorum replication requires that:

    $|Qr| + |Qw| > n$, read and write quorums always intersect

    $2 * |Qw| > n$, any two write quorums always intersect

‣ Read and Write operations take $Qr$ and $Qw$ into account now

‣ Read and Write operations take $Q_r$ and $Q_w$ into account now

# Data replication

‣ Read and Write operations take $Q_r$ and $Q_w$ into account now
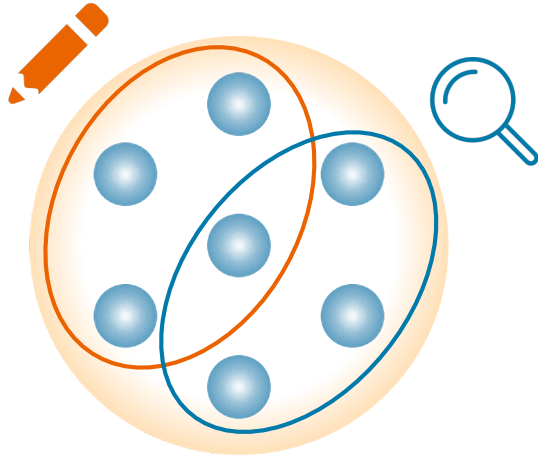
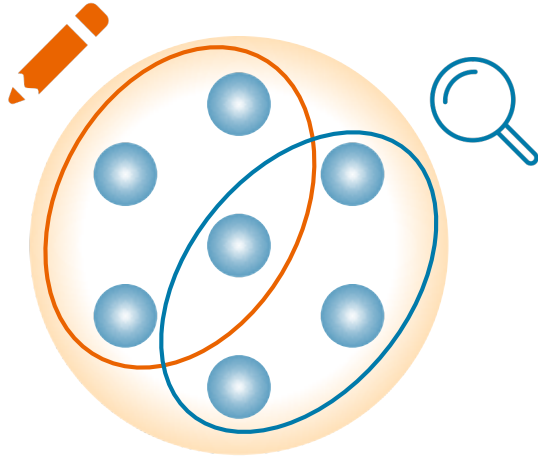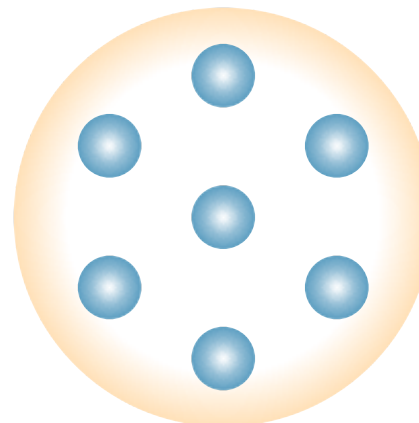‣ Read and Write operations take $Q_r$ and $Q_w$ into account now

# Data replication

## Quorums

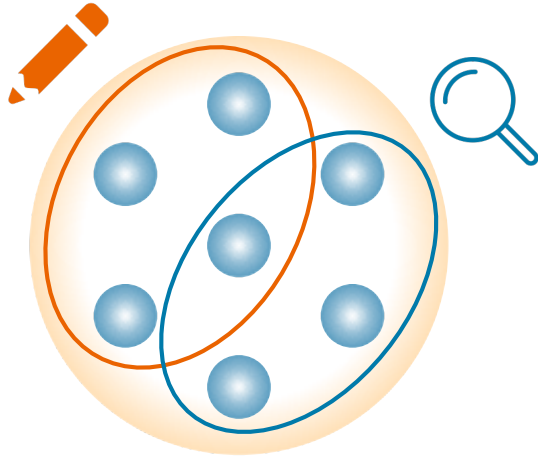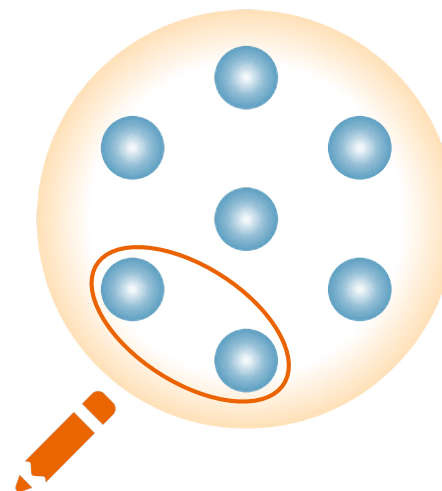‣ Read and Write operations take $Q_r$ and $Q_w$ into account now

▸ Read and Write operations take $Q_r$ and $Q_w$ into account now

‣ Read and Write operations take $Qr$ and $Qw$ into account now



$$|Qr| + |Qw| > n$$

‣ Read and Write operations take $Q_r$ and $Q_w$ into account now



$$|Q_r| + |Q_w| > n$$

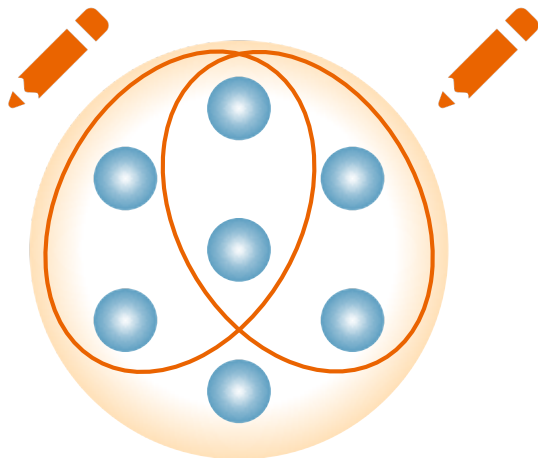‣ Read and Write operations take Qr and Qw into account now



$$|Qr| + |Qw| > n$$

‣ Read and Write operations take Qr and Qw into account now



$$|Qr| + |Qw| > n$$

‣ Read and Write operations take Qr and Qw into account now



$$|Qr| + |Qw| > n$$

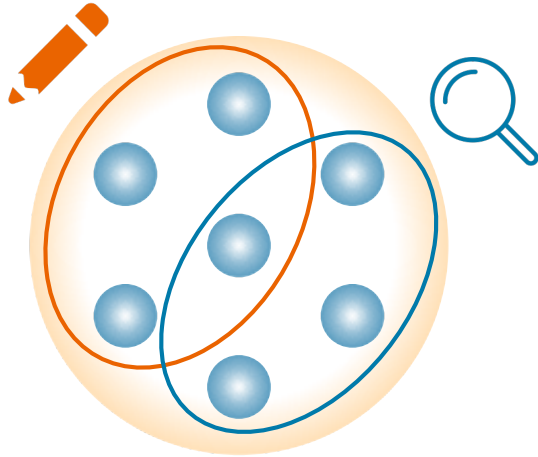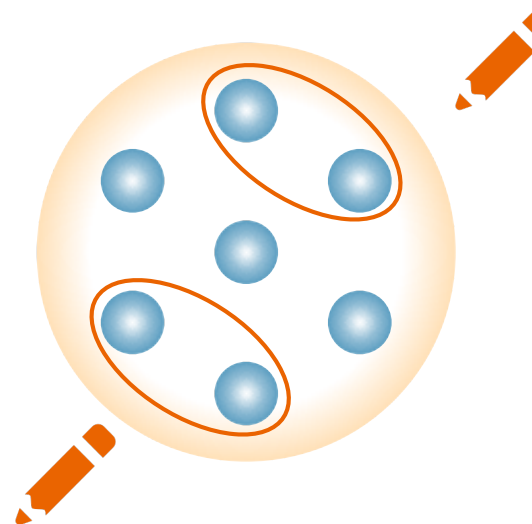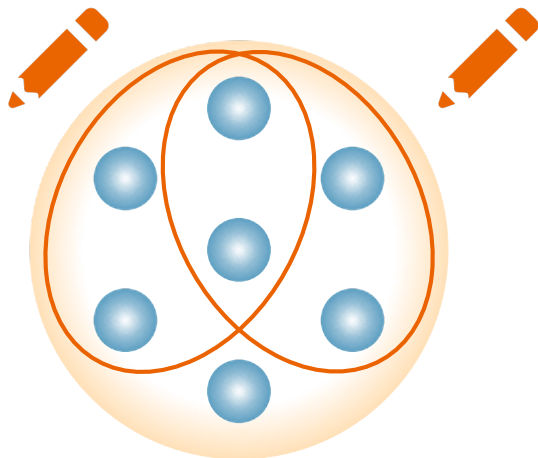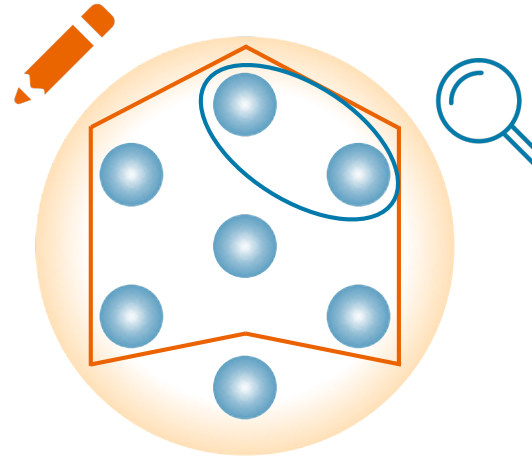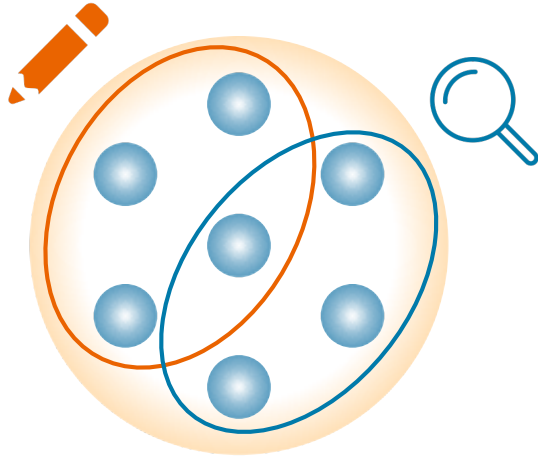‣ Read and Write operations take Qr and Qw into account now
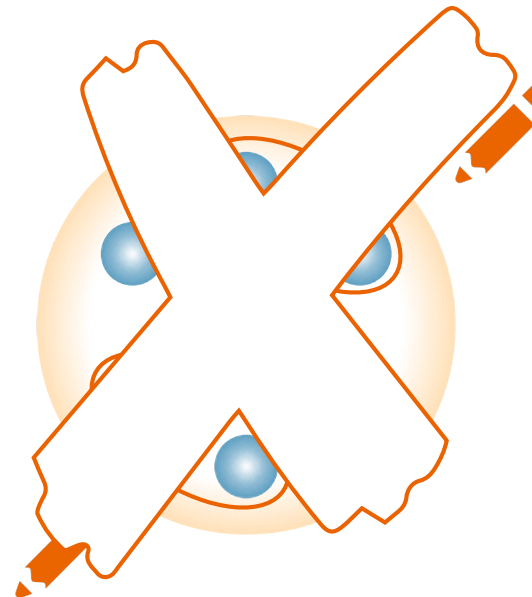


$$|Qr| + |Qw| > n$$

# Data replication

## Quorums

▸ Read and Write operations take Qr and Qw into account now



$$|Qr| + |Qw| > n$$

▸ Read and Write operations take Qr and Qw into account now



$$|Qr| + |Qw| > n$$

$$2 * |Qw| > n$$

‣ Read and Write operations take Qr and Qw into account now



$$|Qr| + |Qw| > n$$

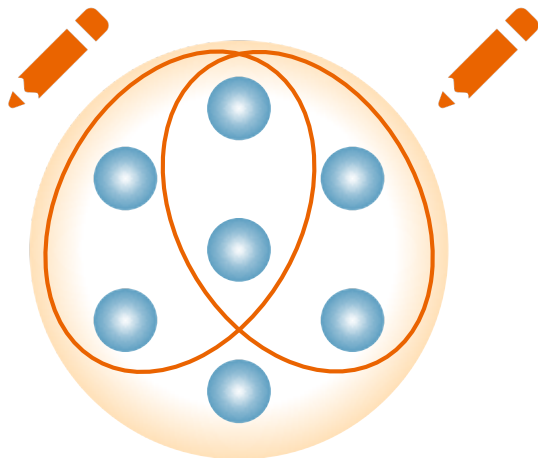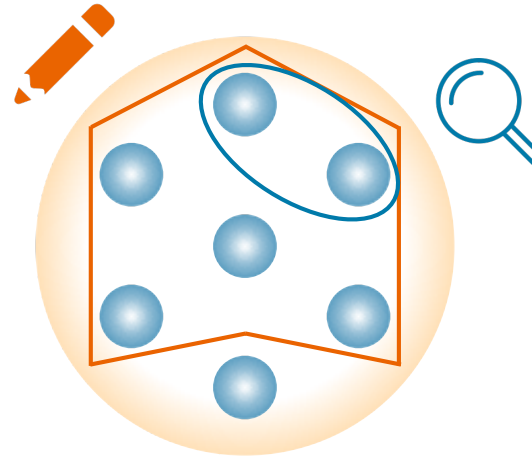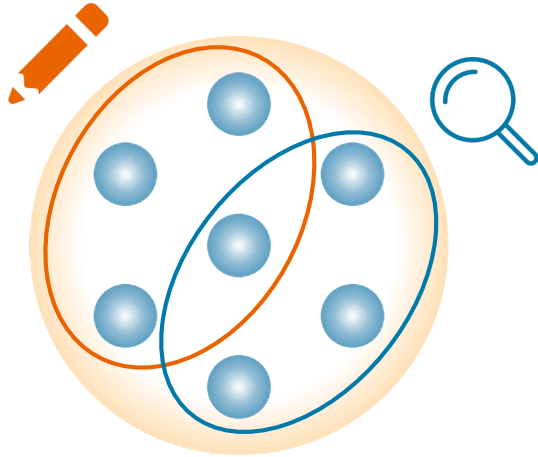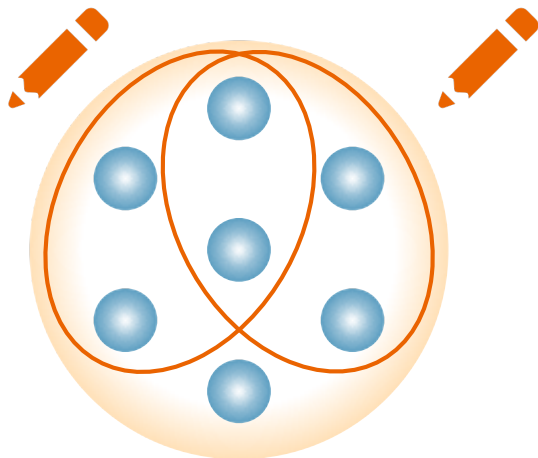$$2 * |Qw| > n$$

‣ Read and Write operations take Qr and Qw into account now

‣ Write operations now need to read from a Qr

At client proxies

read (value)
    SetV = _read (x, v) from a Qr
    vmax = largest (_,v) from SetV
    (value, vmax) from SetV

write (value)
    SetV = _read (_, v) from a Qr
    vmax = largest (_, v) from SetV
    _write (value, vmax+1) to a Qw

At replica i

_read(value, version)
    value = xi
    version = vi

_write (value, version)
    if version > vi then
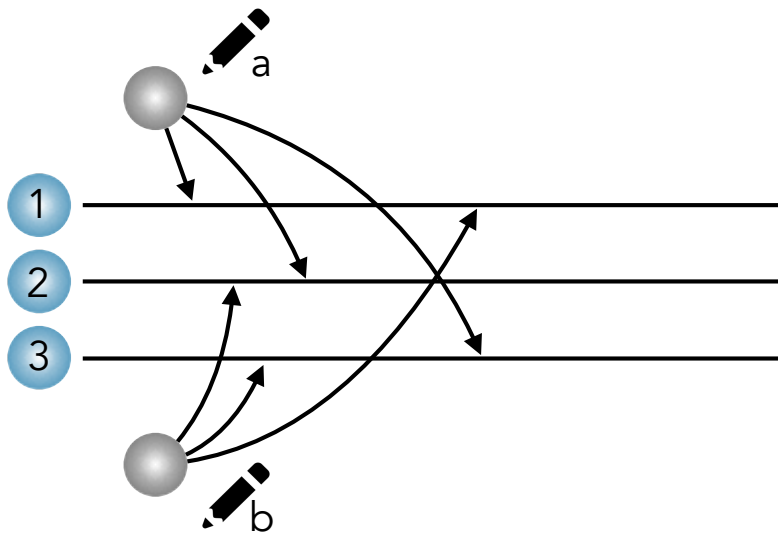        xi = value
        vi = version

# Data replication

‣ The use of quorums allows for trade-offs in several system aspects

‣ With omission faults <u>fault tolerance</u> can be maximised using strict majority quorums:

    ‣ $|Qr| = |Qw| = \lceil (n+1)/2 \rceil$

    ‣ $|Qw| = \lceil (n+1)/2 \rceil$ also leads to the least expensive write operations

‣ $|Qr|$ can determine <u>workload bias</u>; it determines the cost of reads and impacts the cost of writes

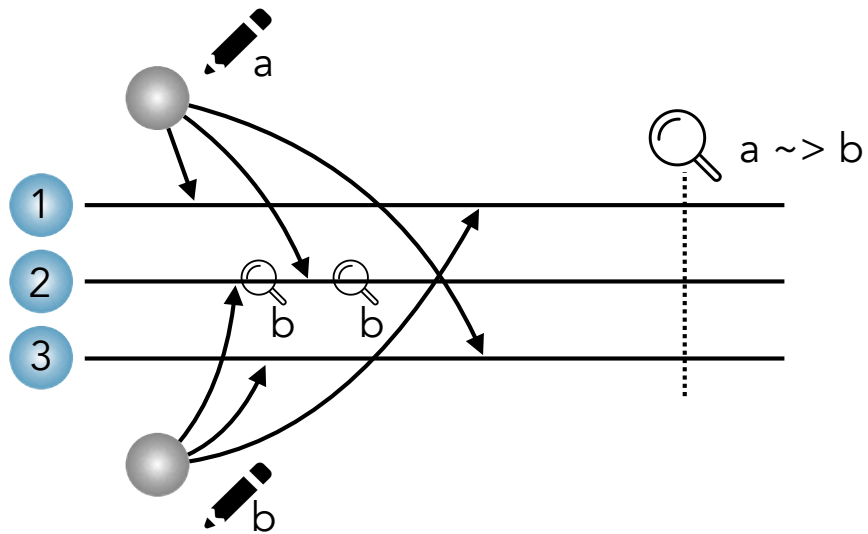‣ For how it may impact throughput, latency, and network load see this chapter's reading material

‣ Monotonic increasing versions + PID's

‣ Monotonic increasing versions + PID's

‣ Monotonic increasing versions + PID's

‣ Monotonic increasing versions + PID's

‣ Monotonic increasing versions + PID's



‣ Globally synchronised clocks

‣ Monotonic increasing versions + PID's



‣ Globally synchronised clocks

‣ Causal ordering + PID's

‣ However, the lack of concurrency control of write operations may lead to basic semantics "inconsistencies"

‣ However, the lack of concurrency control of write operations may lead to basic semantics "inconsistencies"

‣ Reading **within** a concurrent write may lead to **unexplainable** results

‣ However, the lack of concurrency control of write operations may lead to basic semantics "inconsistencies"

    ‣ Reading **within** a concurrent write may lead to **unexplainable** results



‣ Consider an application that manages a non-decreasing variable

# Data replication

‣ However, the lack of concurrency control of write operations may lead to basic semantics "inconsistencies"

  ‣ Reading **within** a concurrent write may lead to **unexplainable** results

‣ Consider an application that manages a non-decreasing variable
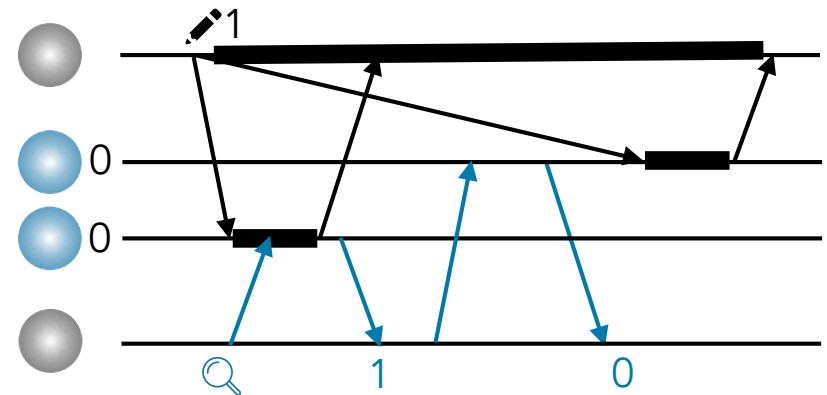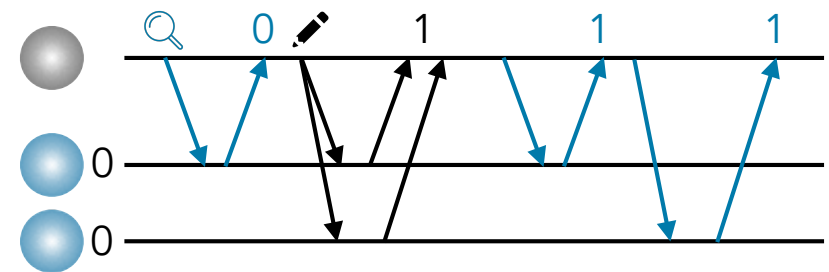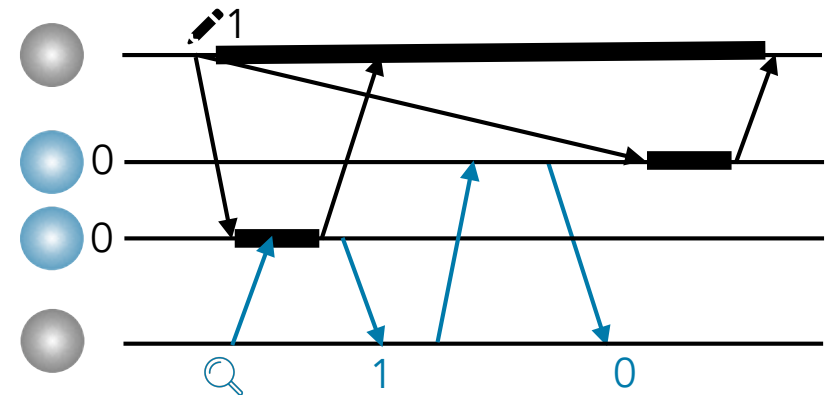
  ‣ One process (sequential) execution is OK
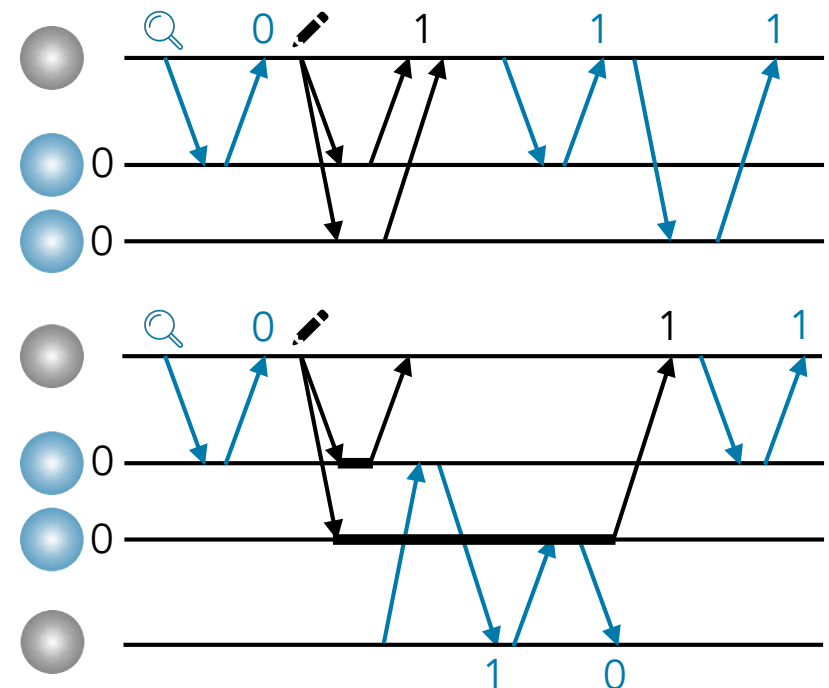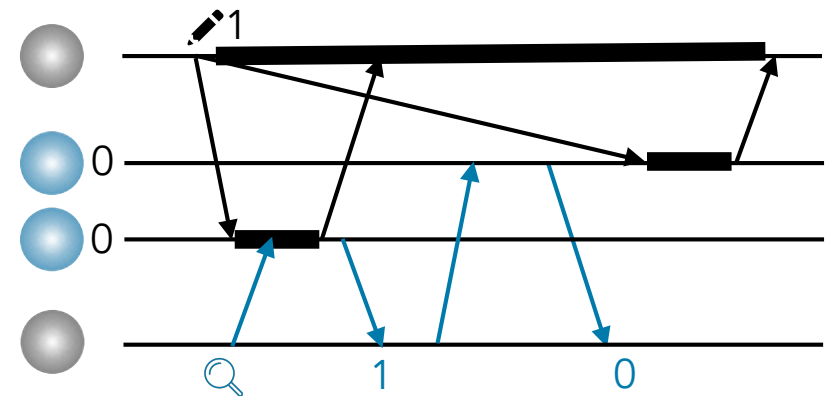
‣ However, the lack of concurrency control of write operations may lead to basic semantics "inconsistencies"

  ‣ Reading **within** a concurrent write may lead to **unexplainable** results



‣ Consider an application that manages a non-decreasing variable

  ‣ One process (sequential) execution is OK

  ‣ A concurrent execution may easily violate the application semantics

# Data replication

‣ Read and Write operations take Qr and Qw into account now

### At client proxies

```
read (value)
    SetL = _getlock() from a Qr
    if |SetL| == |Qr|
        SetV = _read (x, v) from a Qr
        vmax = largest (_,v) from SetV
        (value, vmax) from SetV
        res = OK
    else
        res = error
    _freelock() from SetL
    return res


write (value)
    SetL = _getlock() from a Qw
    if |SetL| == |Qw|
        SetV = _read (_, v) from a Qr
        vmax = largest (_, v) from SetV
        _write (value, vmax+1) to a Qw
        res = OK
    else
        res = error
    _freelock() from SetL
    return res
```

### At replica i

```
_getlock()
    if lockedi == False
        lockedi = True
        return i
    else
        return error

_freelock()
    lockedi = False

_read(value, version)
    value = xi
    version = vi

_write (value, version)
    if version > vi then
        xi = value
        vi = version
```

# Models of distributed systems and related faults

▸ D. Gifford

"Weighted voting for replicated data"

SYMPOSIUM ON OPERATING SYSTEMS PRINCIPLES, 1979

▸ M. Whittaker, A. Charapko, J. Hellerstein, H. Howard, I. Stoica

"Read-Write Quorum Systems Made Practical"

PAPOC'21, 2021