# Sistemas Distribuídos

José Orlando Pereira

Departamento de Informática
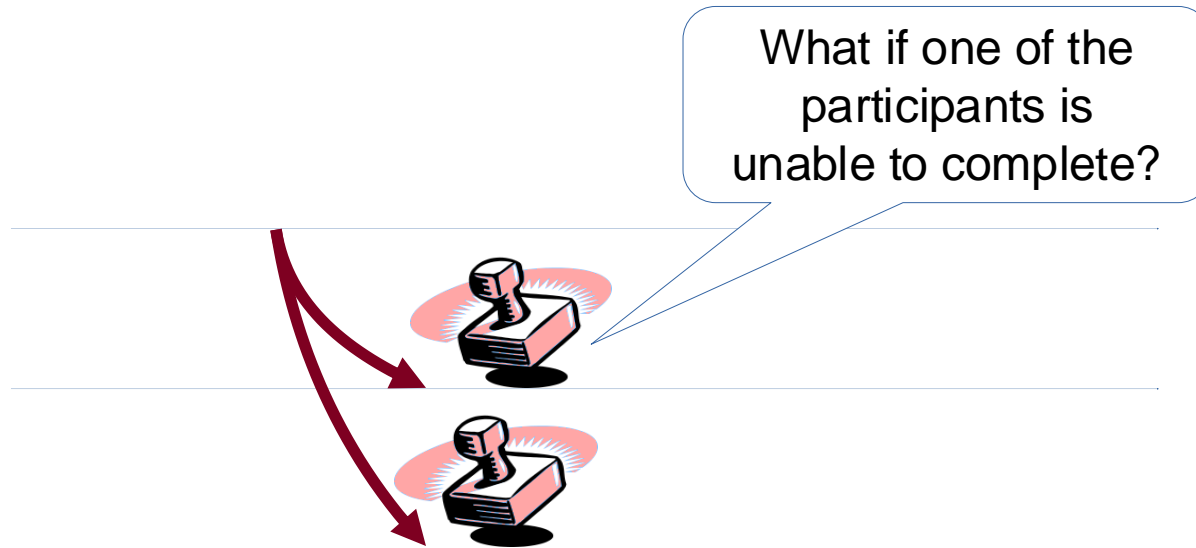Universidade do Minho
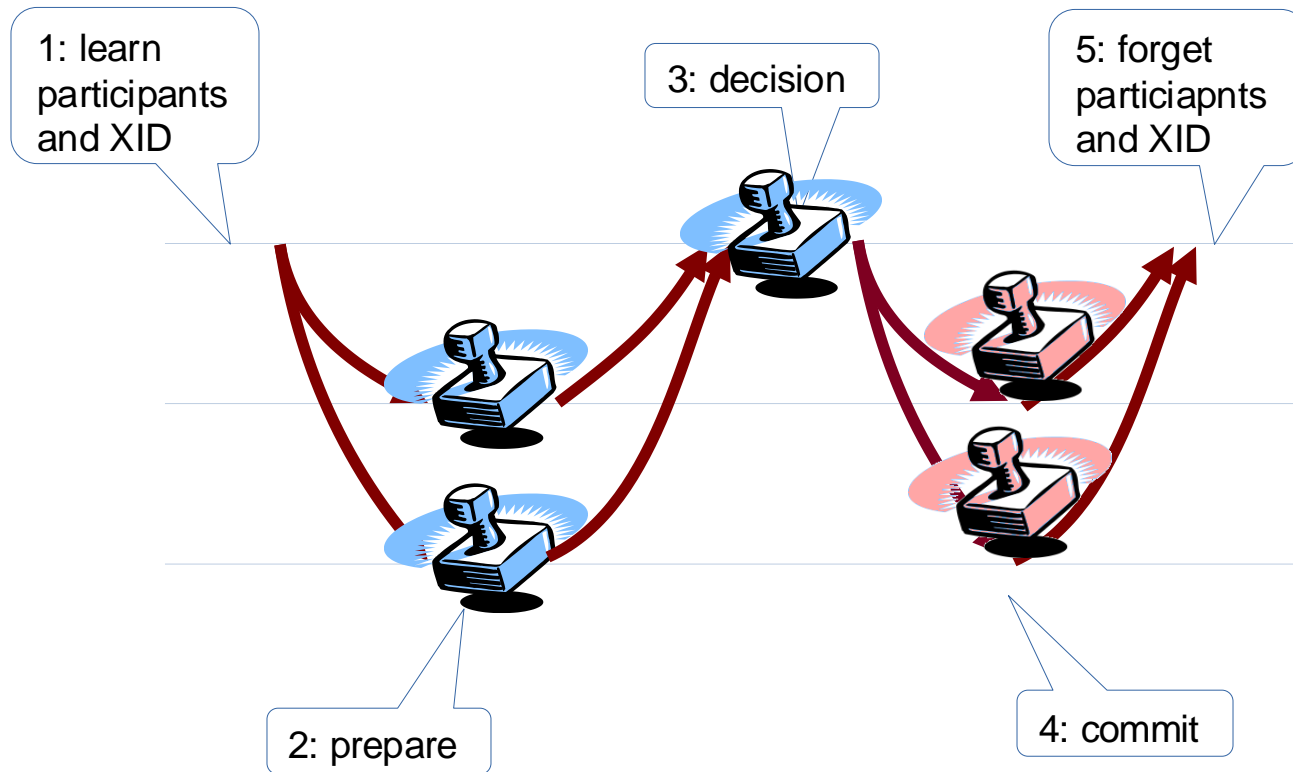
# Fault tolerance

- A distributed system is composed by autonomous computing elements and can tolerate <u>faults</u> to avoid <u>failure</u>

- Fault model describes types of faults:

    - Omissive: process crash, lost message, …

    - Assertive (a.k.a. Byzantine): corrupted messages, …

- Fault model describes the number of faults:

    - Example: Number of processes that can crash
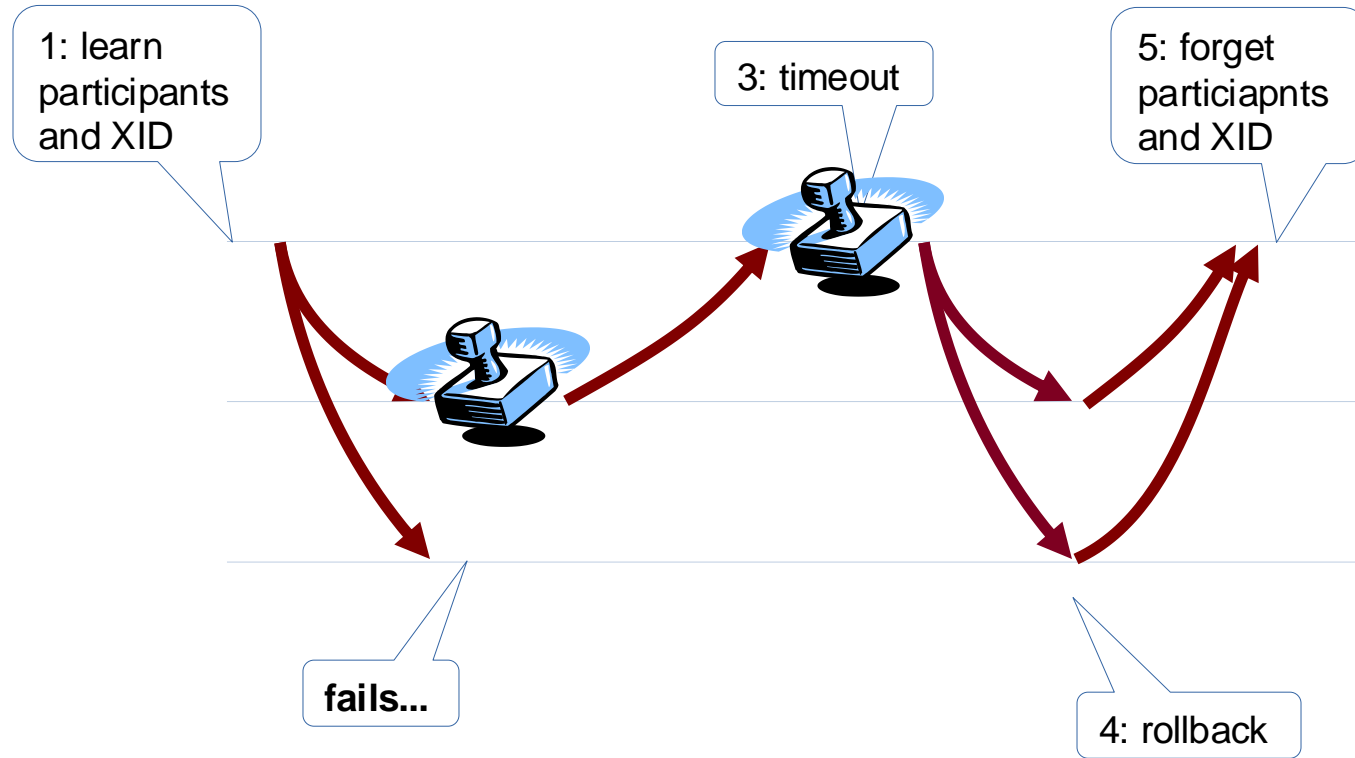
# Transactional commit

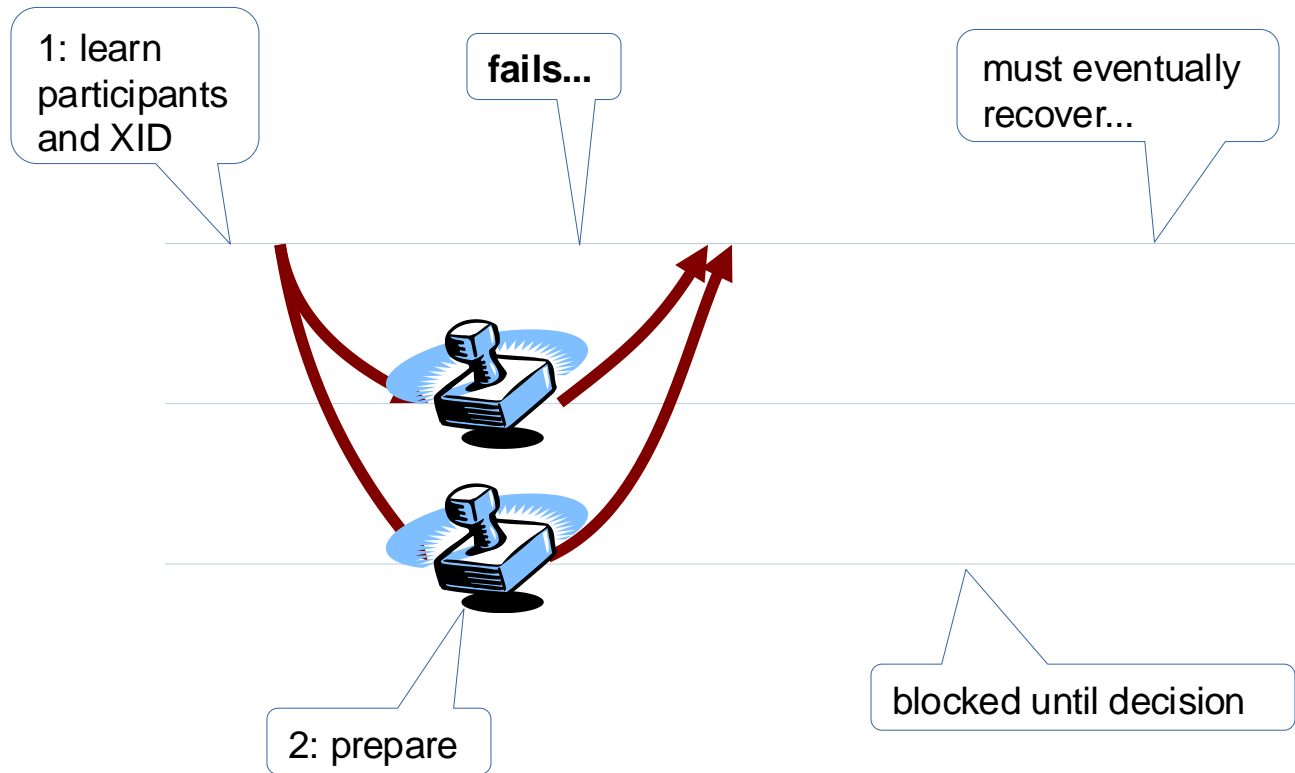- Coordinate multiple irreversible actions across a distributed system:



What if one of the participants is unable to complete?

# 2-phase commit (2PC)



1: learn participants and XID

2: prepare

3: decision

4: commit

5: forget participiapnts and XID

# 2PC: Participant failure



1: learn participants and XID

3: timeout

5: forget particiapnts and XID

fails...

4: rollback

# 2PC: Coordinator failure

1: learn participants and XID

fails...

must eventually recover...

2: prepare

blocked until decision

# 2PC in systems

Coordinator

Resource 1

Resource 2

Client

begin()

Business logic

2PC

phase 1  phase 2

commit()

# Transactional RPC

# Transactional RPC



add resource 2

add resource 1

phase 1 phase 2

Manager

Resource 1

xid+resource 2

Resource 2

xid+resource 1

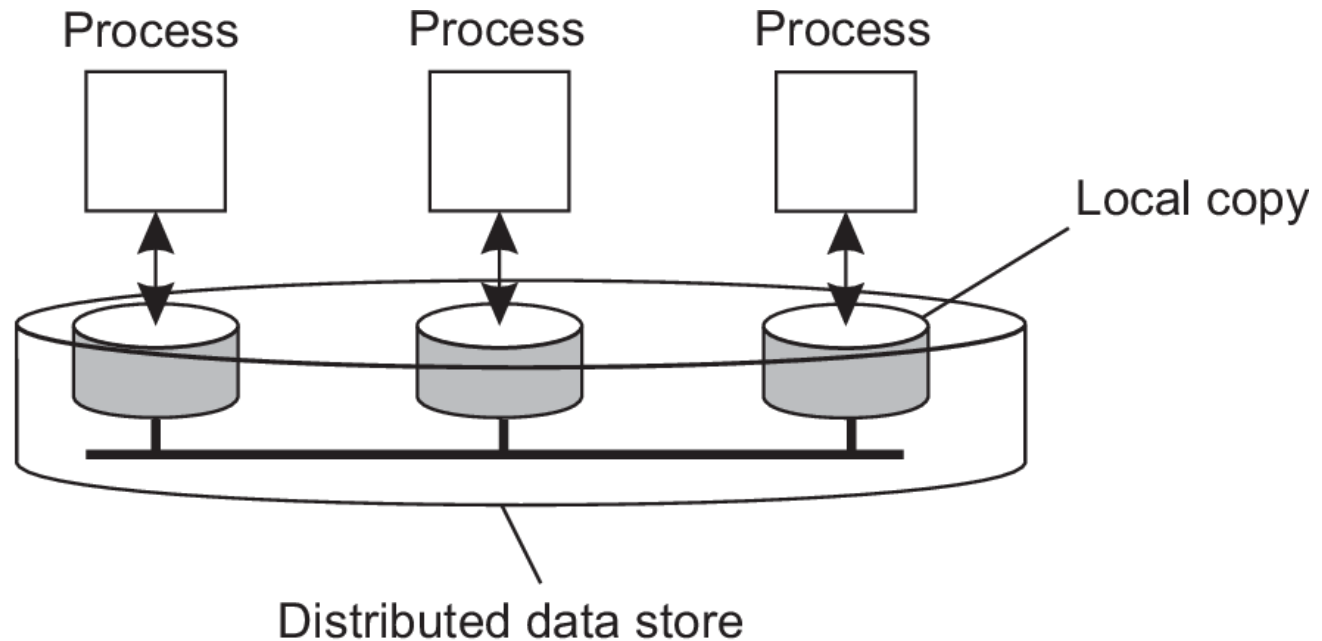Nested invocation

Client

begin()

commit()

# Summary

- Distributed transactions with 2-phase commit (2PC) support agreement in systems with faults

  - Limited to crash-recovery of the coordinator

- Is widely used in enterprise middleware for application integration

# Replication

- Keep multiple copies of the same data or service
  - Distribute the load for scalability
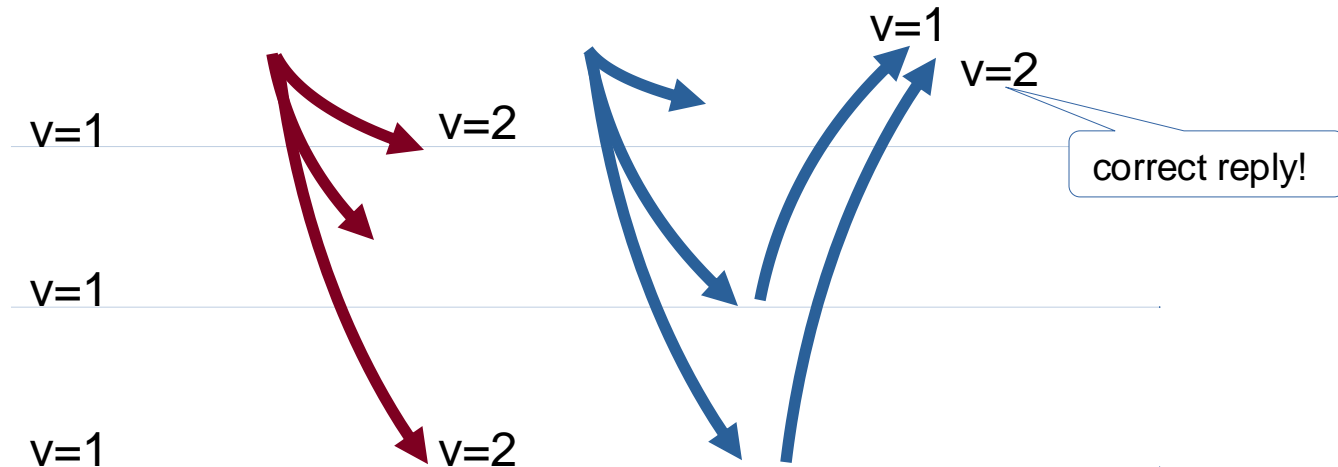  - Tolerate server faults

# Replication

- Naive solution: write then propagate
  - state may diverge
  - clients observe paradoxes when reading
  - not fault-tolerant
- 2PC: Correct, but progress only with all up (tolerates reboots, not crashes)
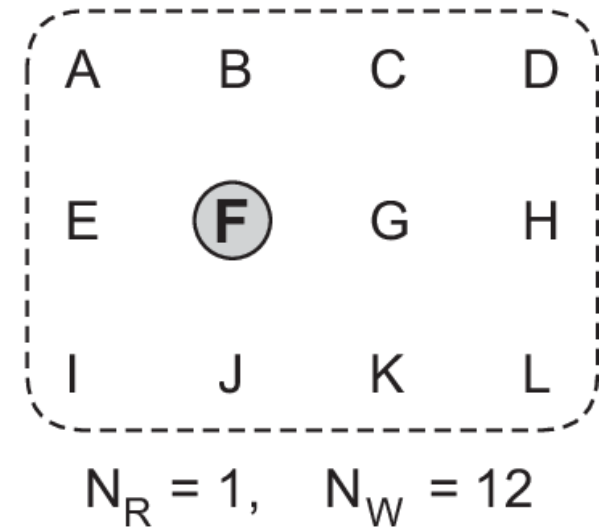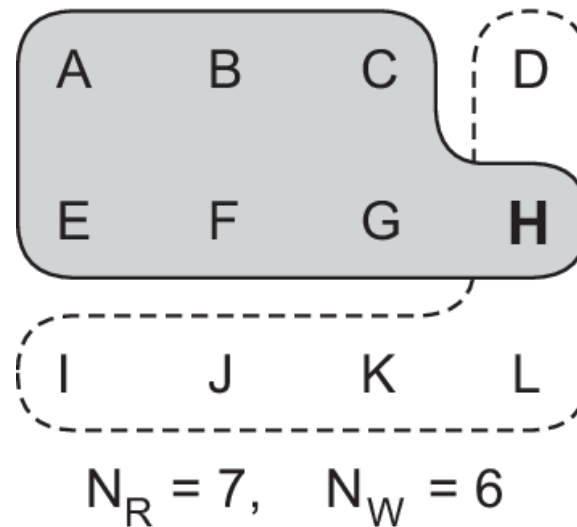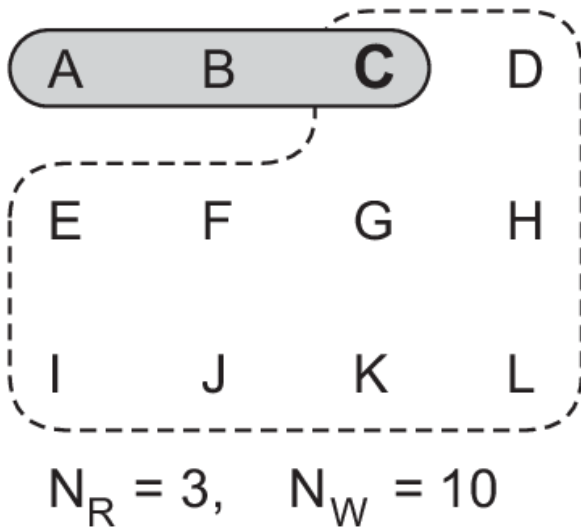
# Replication

- Assume that:
    - operations are reads and writes
    - we keep a timestamp with each item
- It might be possible to read and write from fewer processes...

# Quorum

- Assume 2-phase protocol for writing (phase 1 == read)
- Quorum rules for replicated data:
  - $N_R + N_W > N \rightarrow$ readers get the latest value
  - $N_W > N/2 \rightarrow$ concurrent writers conflict



$N_R = 3, \quad N_W = 10$     $N_R = 7, \quad N_W = 6$     $N_R = 1, \quad N_W = 12$

# Quorum

- Additional rules for fault-tolerance when assuming at most *f* faults:

    - $N_R + f \leq N \rightarrow$ readers never block

    - $N_W + f \leq N \rightarrow$ writers never block

- Can be configured to ensure both or either of them

- Typical solution is having a majority:

    - *$N_R = N_W = f + 1$*

    - *$N = 2f + 1$*

- Examples: N=3 for f=1, N=5 for f=2, ….

# Summary

- Flexible and efficient solution for data replication
  - Example: Amazon Aurora DB
- Wasteful when there are multiple concurrent writers:
  - At most one of multiple write operation can be accepted
  - But it can happen that none is accepted if each operation is applied in less than $N_W$ servers