

CS 347:
Distributed Databases and
Transaction Processing
**Notes07: Reliable
Distributed
Database Management**

Hector Garcia-Molina

Reliable distributed database management

- Reliability
- Failure models
- Scenarios

Reliability

- Correctness
 - Serializability
 - Atomicity
 - Persistence
- Availability

Types of failures

- Processor failures
 - Halt, delay, restart, bezerk, ...
- Storage failures
 - Volatile, non-volatile, atomic write, transient errors, spontaneous failures
- Network failures
 - Lost message, out-of-order messages, partitions, bounded delay

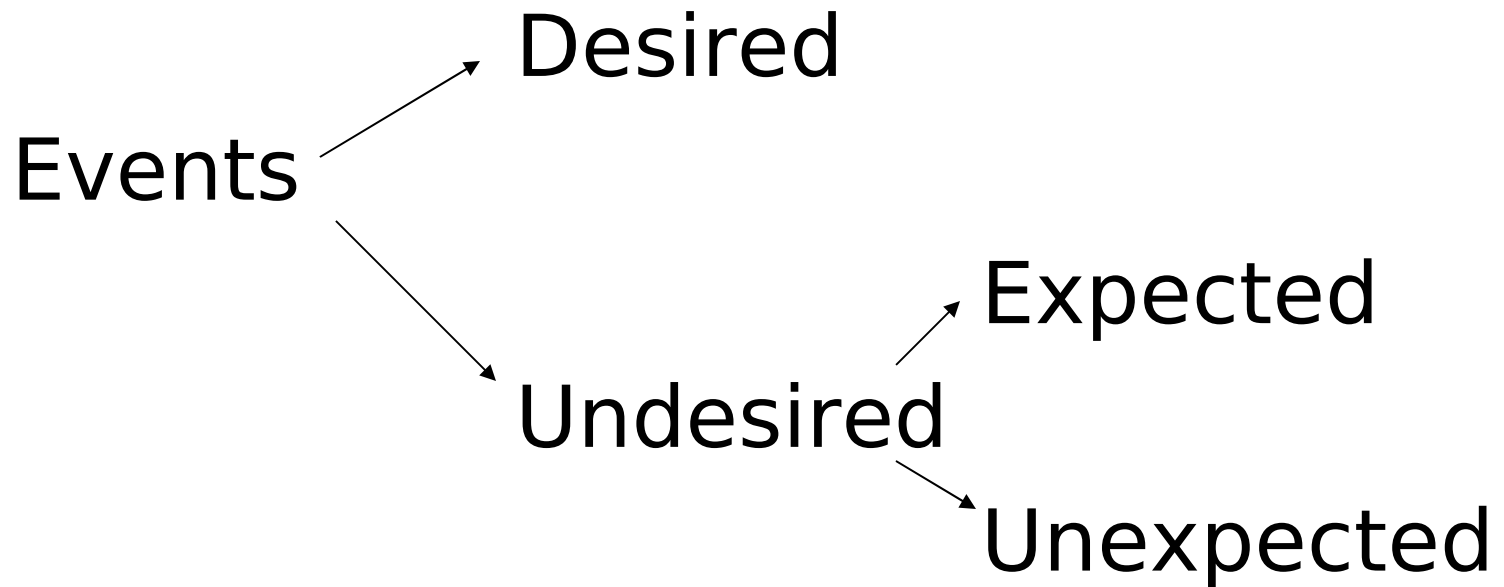
More Types of failures

- Malevolent failures
- Multiple failures
- Detectable failures

Failure models

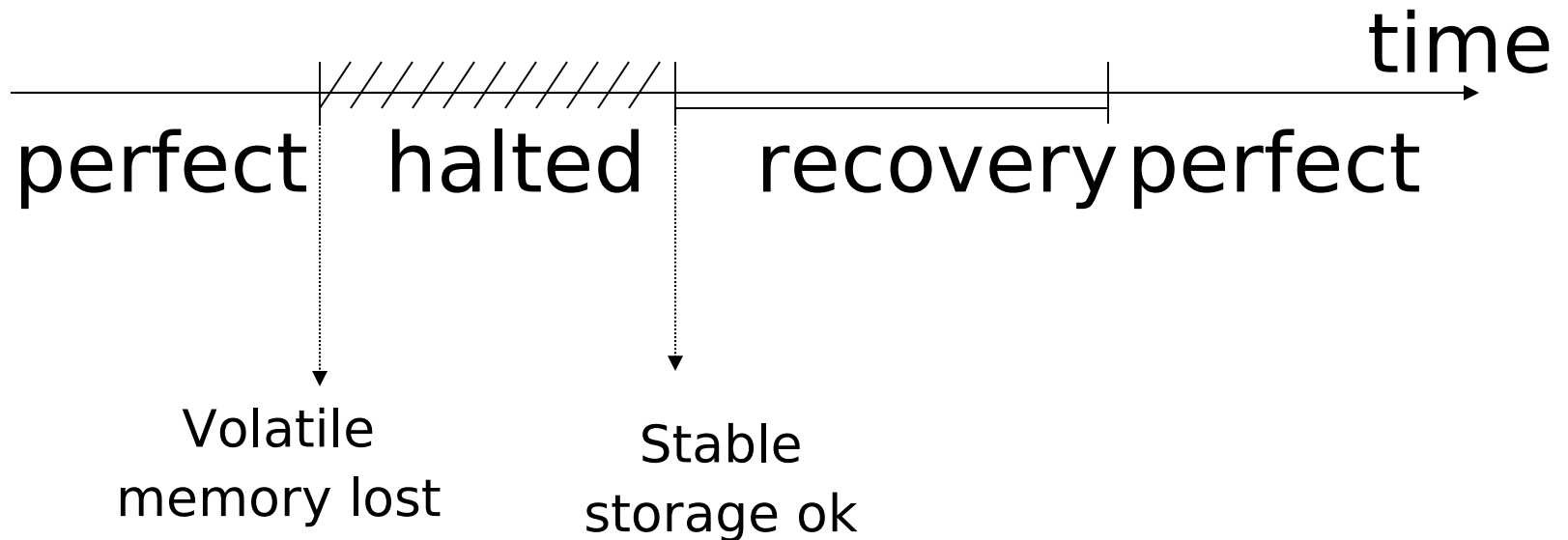
- Cannot protect against everything
- Unlikely failures (e.g., flooding in the Sahara)
- Expensive to protect failures
(e.g., earthquake)
- Failures we know how to protect
against
(e.g., message sequence numbers; stable storage)

Failure model:



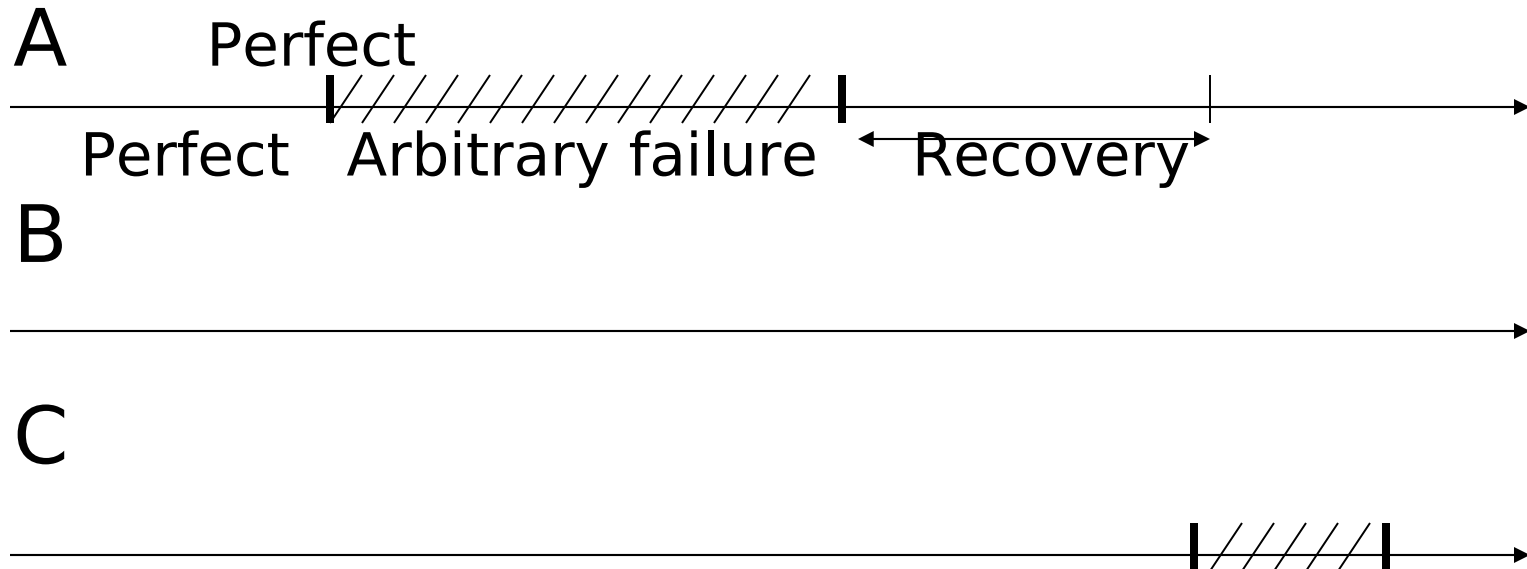
Node models

(1) Fail-stop nodes



Node models

(2) Byzantine nodes



At any given time, at most some fraction f of nodes failed (typically $f < 1/2$ or $f < 1/3$)

Network models

(1) Reliable network

- in order messages
- no spontaneous messages
- timeout T_D

If no ack in T_D sec.



Destination down
(not paused)

I.e., no lost messages, except for node failures

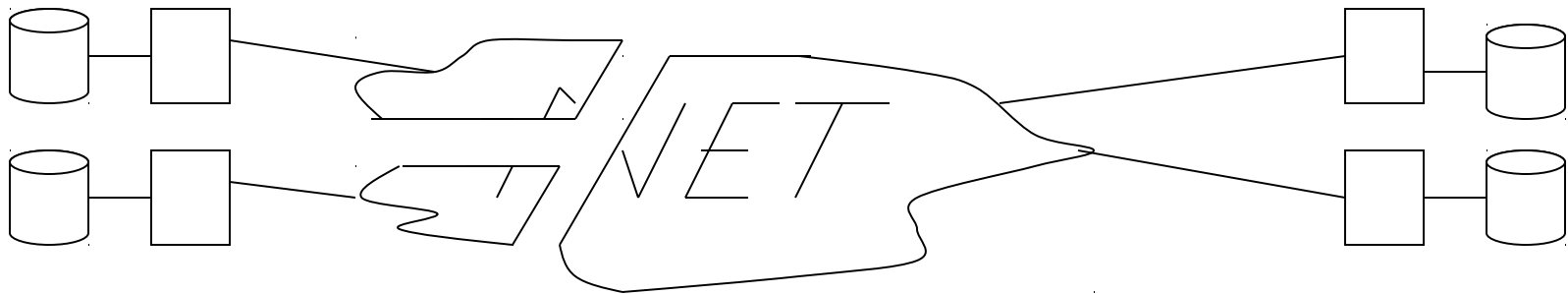
Variation of reliable net:

- Persistent messages
 - If destination down, net will eventually deliver message
 - Simplifies node recovery, but leads to inefficiencies (hides too much)
 - Not considered here

Network models

(2) Partitionable network

- In order messages
- No spontaneous messages



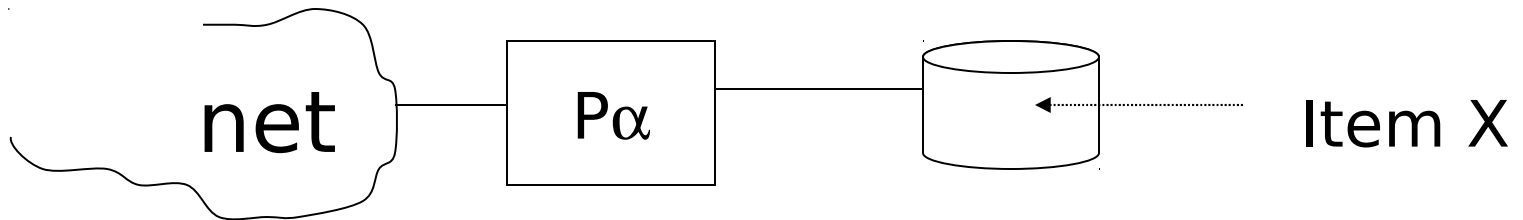
- no timeout; nodes can have different view of failures

Scenarios

- Reliable network
 - Fail-stop nodes
 - No data replication (1)
 - Data replication (2)
- Partitionable network
 - Fail-stop nodes (3)

No Data Replication

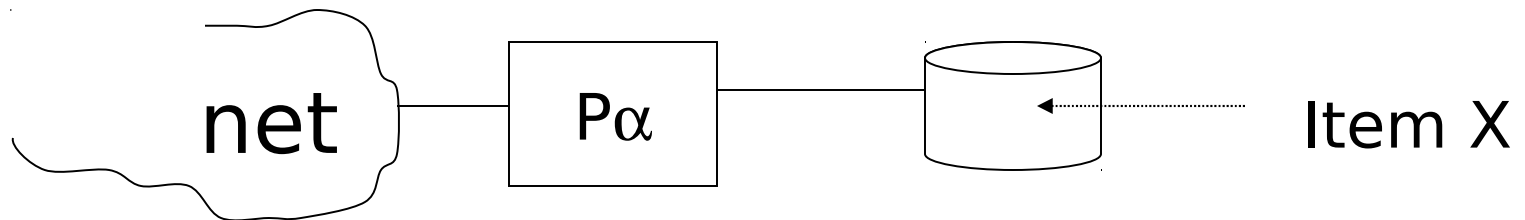
- Reliable network, fail-stop nodes



- Basic idea: node P_α controls X

No Data Replication

- Reliable network, fail-stop nodes

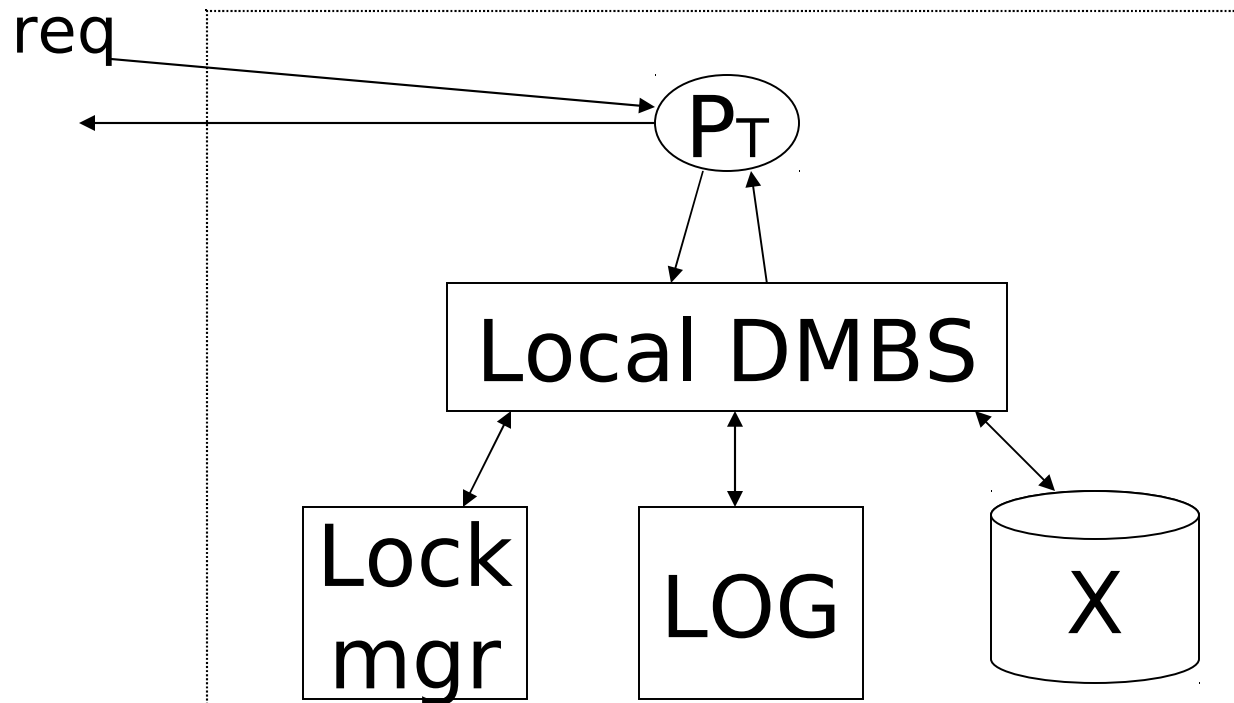


- Basic idea: node P_α controls X
 - Single control point simplifies concurrency control, recovery
 - Not an availability hit:
if P_α down, X unavailable too!

“ P_α controls X ” means

- P_α does concurrency control for X
- P_α does recovery for X

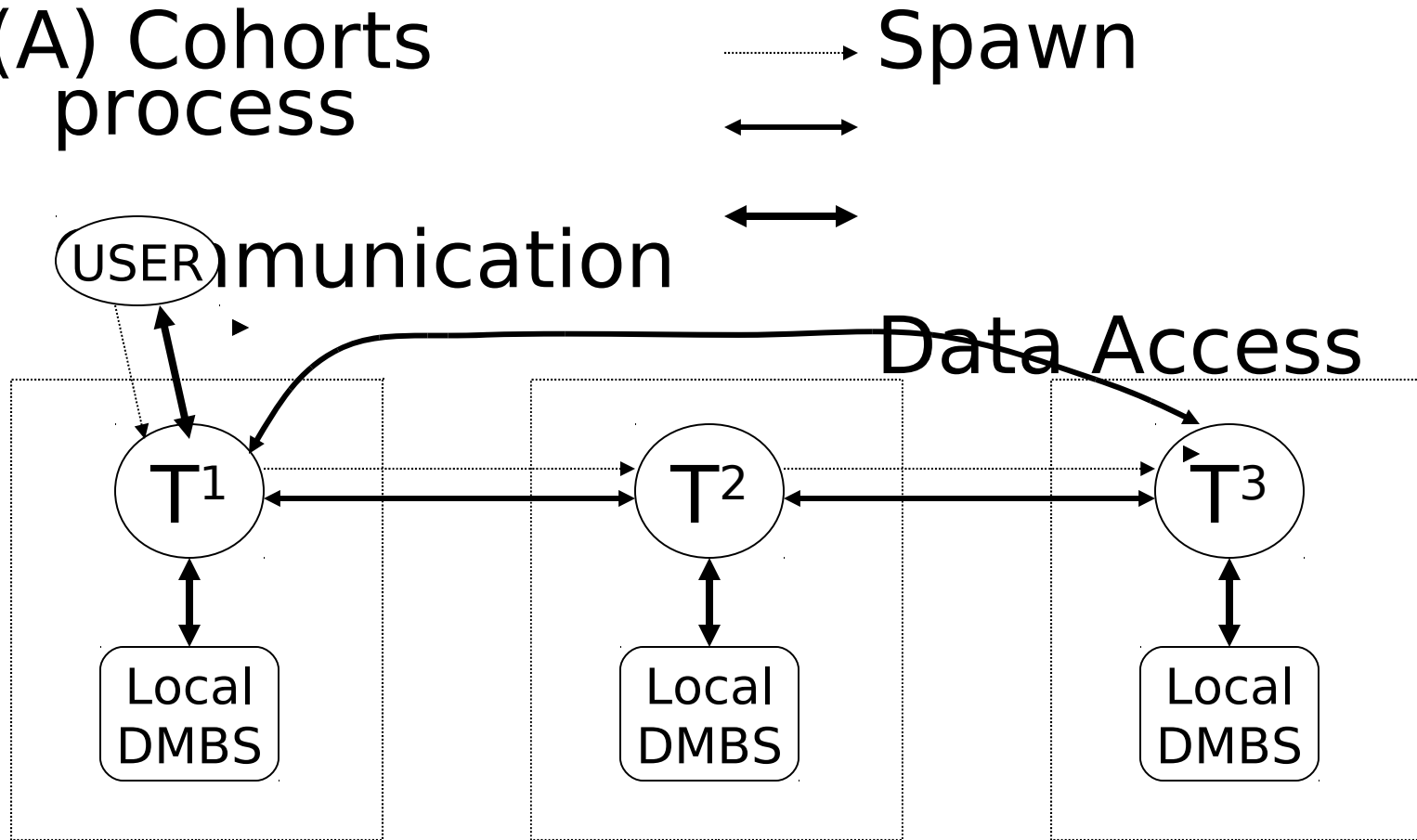
Say transaction T wants to access X:



P_T is process that represents T at this node

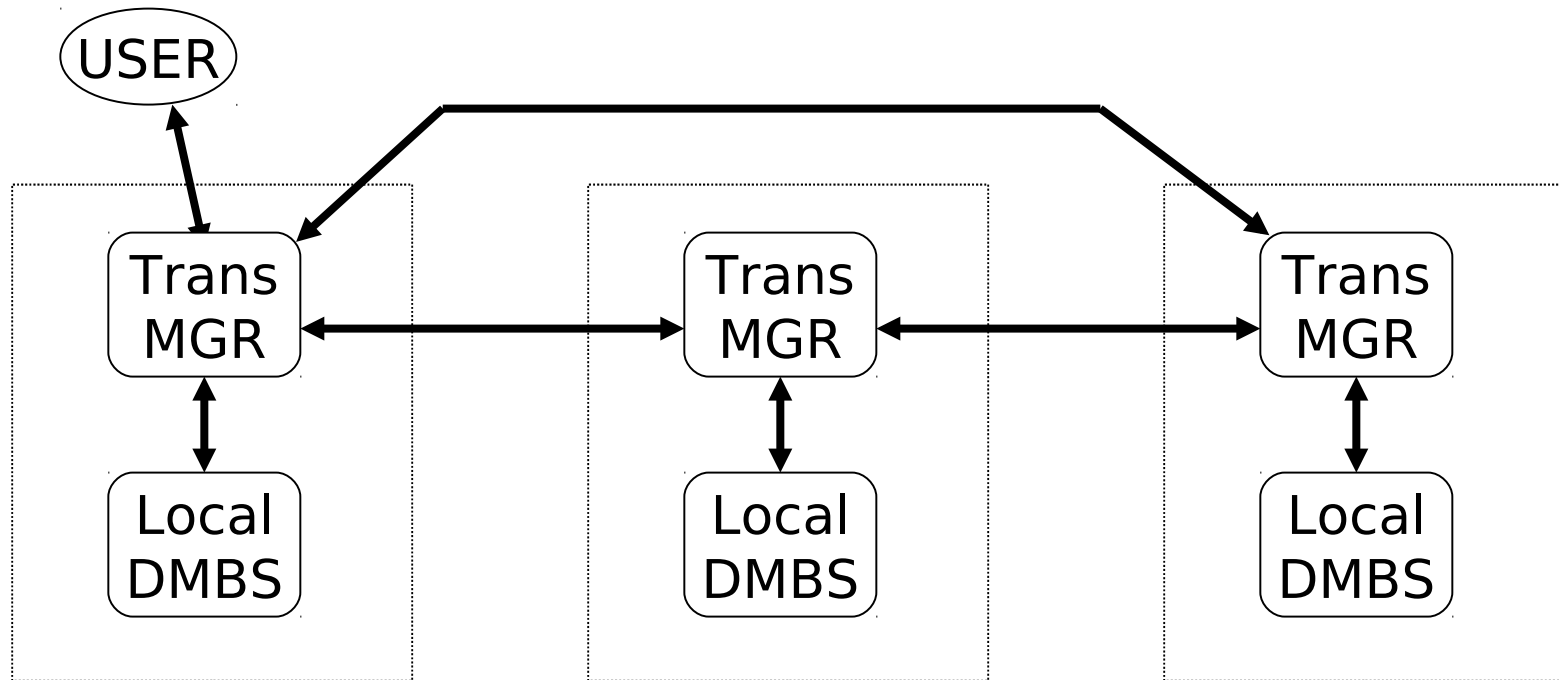
Process models

(A) Cohorts process



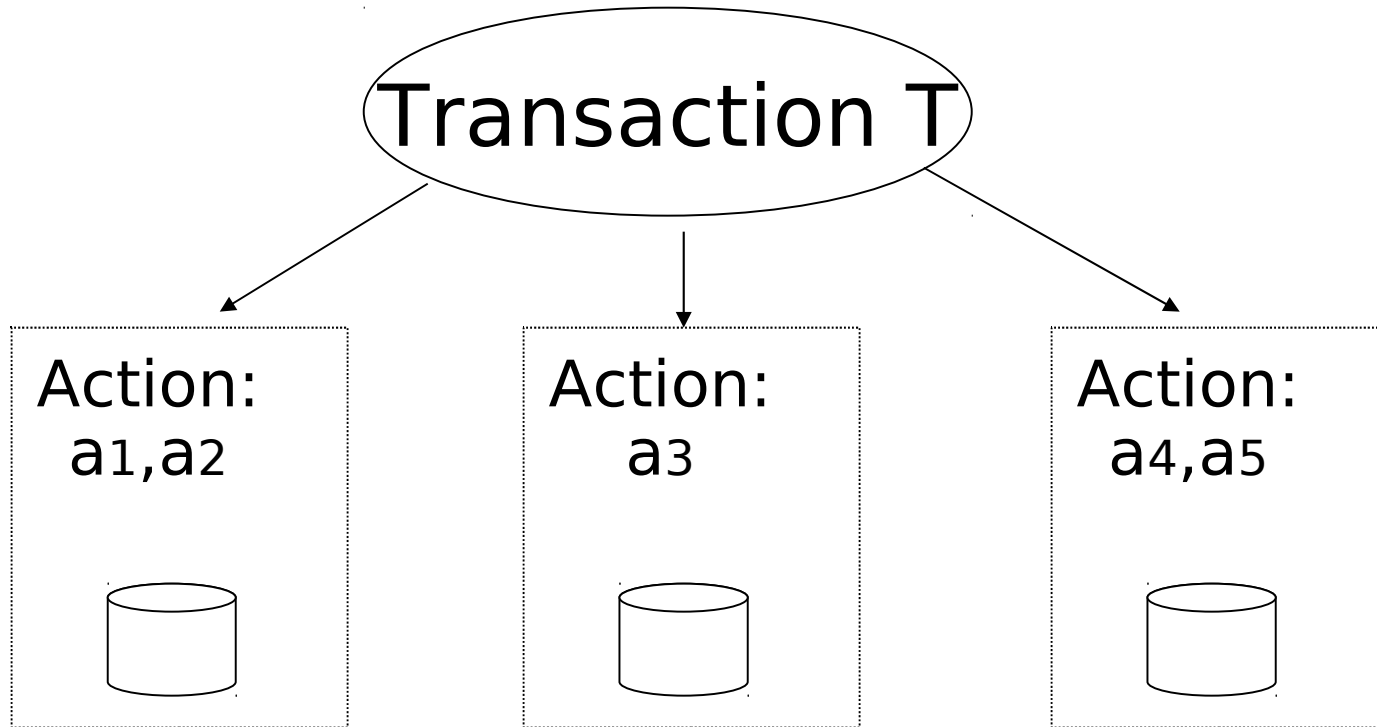
Process models

(B) Transaction servers (manager)



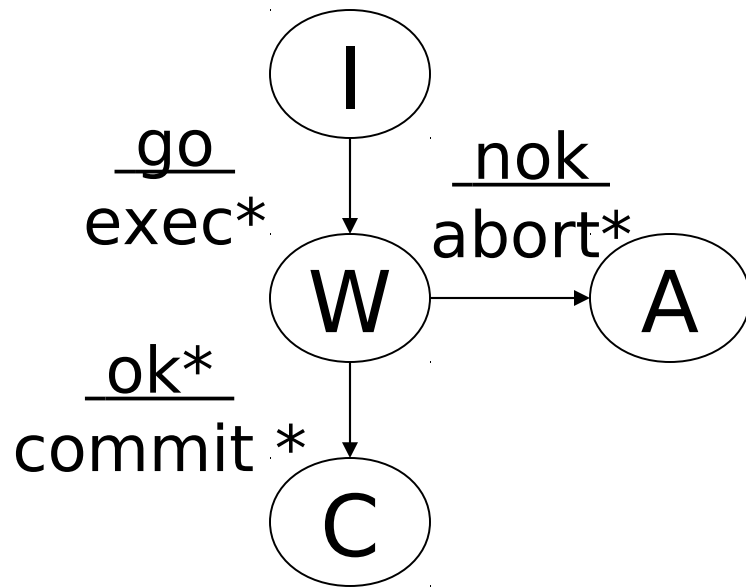
- Cohorts: application code responsible for remote access
- Transaction manager: “system” handles distribution, remote access

Distributed commit problem

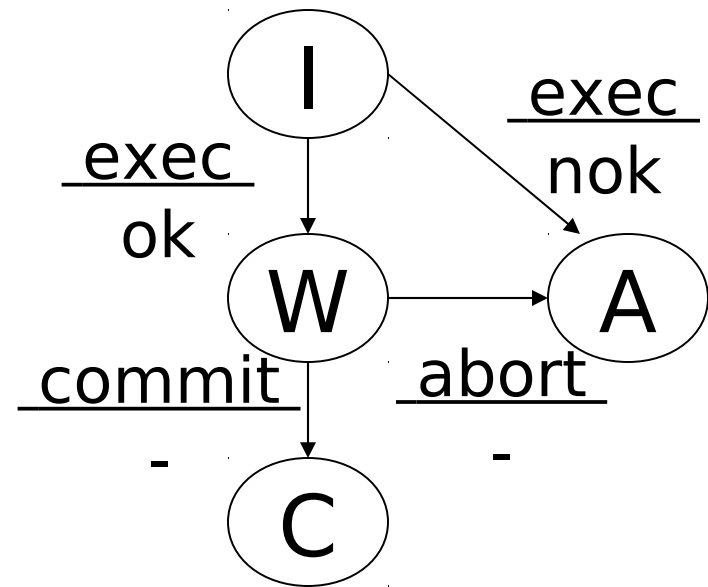


Centralized two-phase commit

Coordinator



Participant



- Notation: Incoming message
Outgoing message
(* = everyone)
- When participant enters “W” state:
 - it must have acquired all resources
 - it can only abort or commit if so instructed by a coordinator
- Coordinator only enters “C” state if all participants are in “W”, i.e., it is certain that all will eventually commit

Handling node failures

- Coordinator and participant logs are used to reconstruct state before failure

-> Example: after participant fails:

Log:

	T1 X undo/recovery info	...	T1 Y info	...	T1 "W" state	
--	----------------------------------	-----	-----------------	-----	--------------------	--

☞ At recovery:

- T₁ is in “W” state
- Obtain X,Y write locks (no read locks!)
- Wait for message from coordinator
(or ask coordinator for outcome)

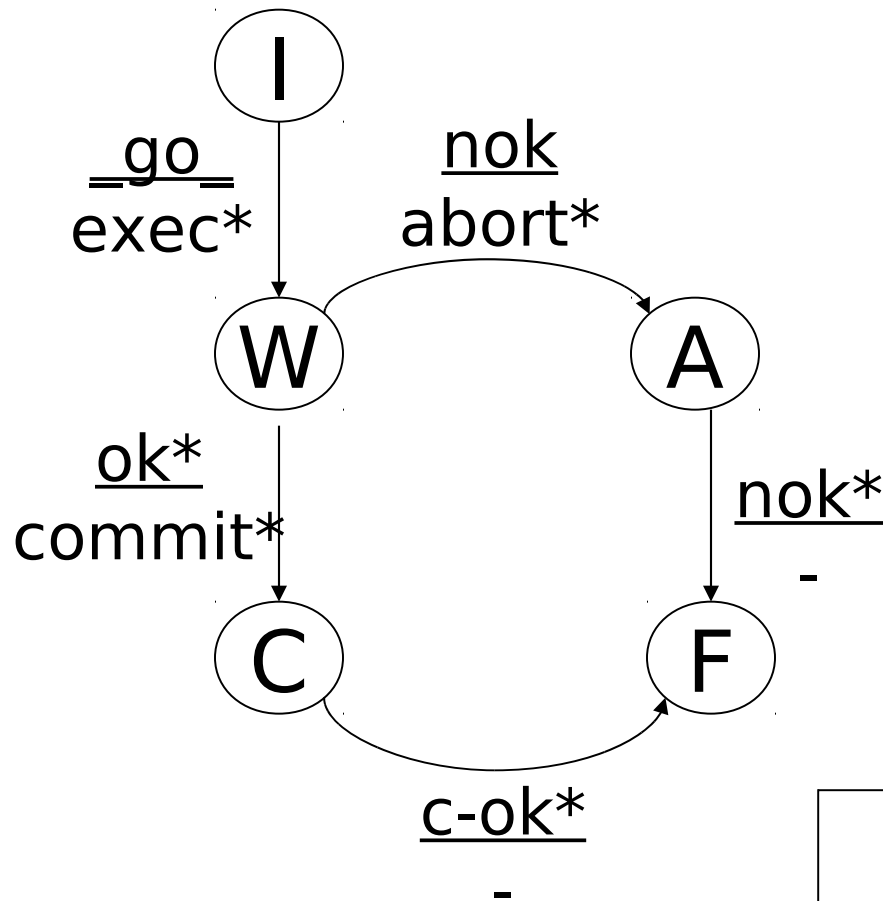
☞ Other examples:

- No “W” record on log \Rightarrow abort T_1
- See “C” record on log \Rightarrow finish T_1

Next

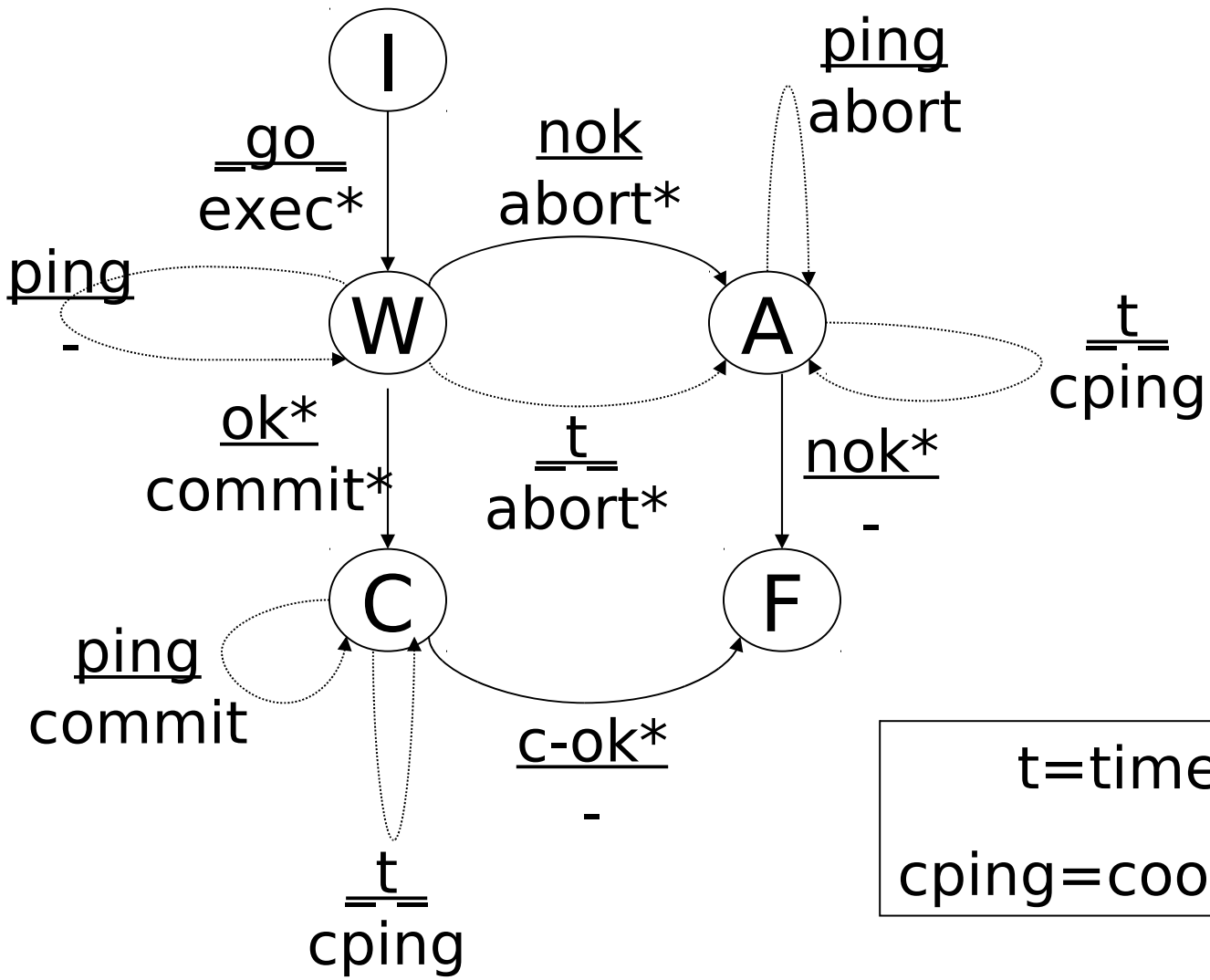
- Add timeouts to cope with messages lost during crashes
- Add finish (“F”) state for coordinator – all done, can forget outcome

Coordinator

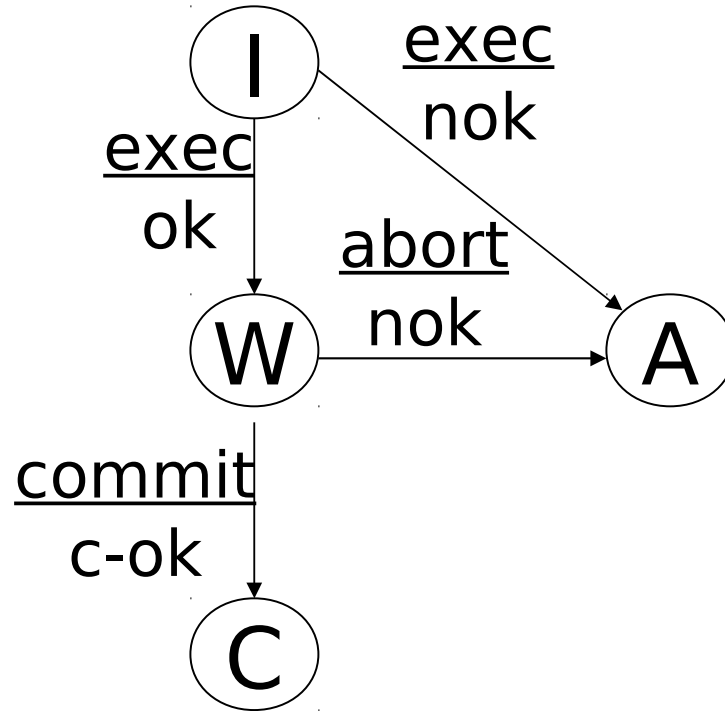


t=timeout
cping=coord. ping

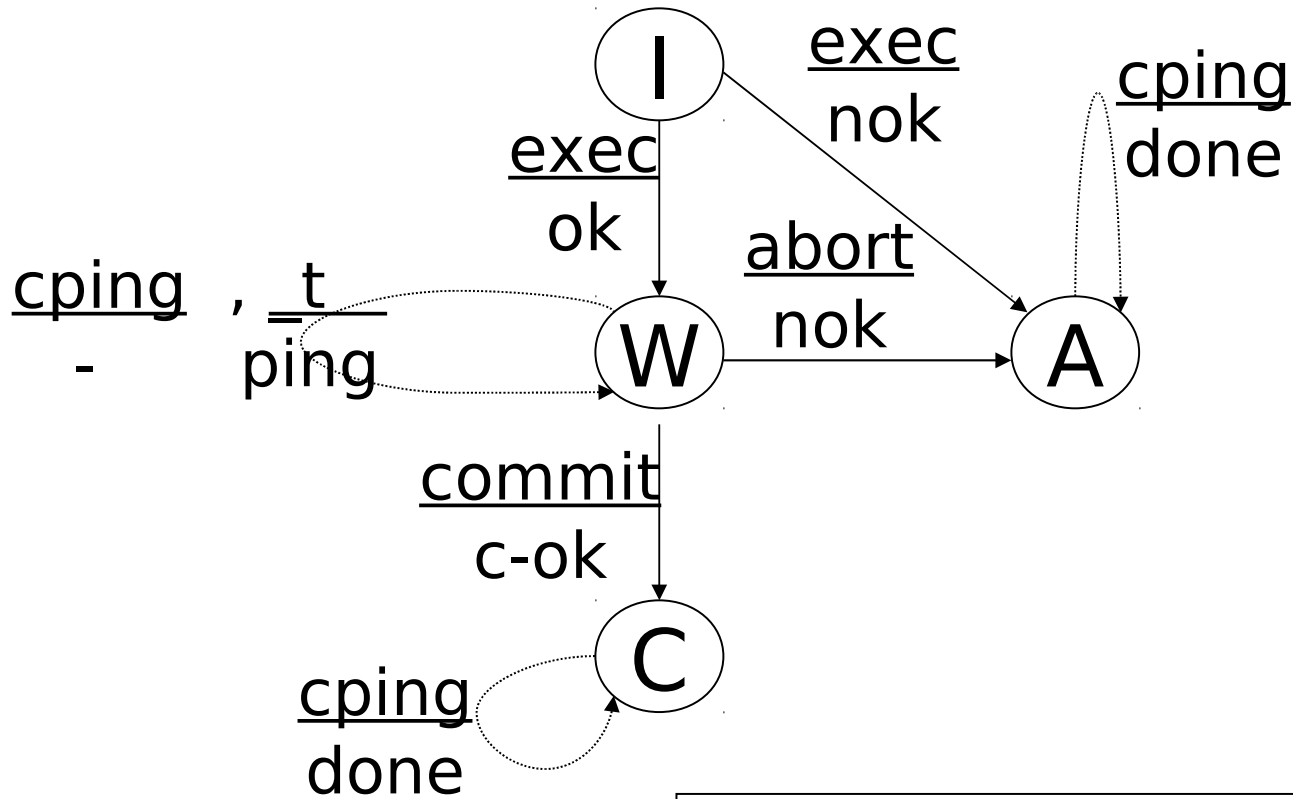
Coordinator



Participant

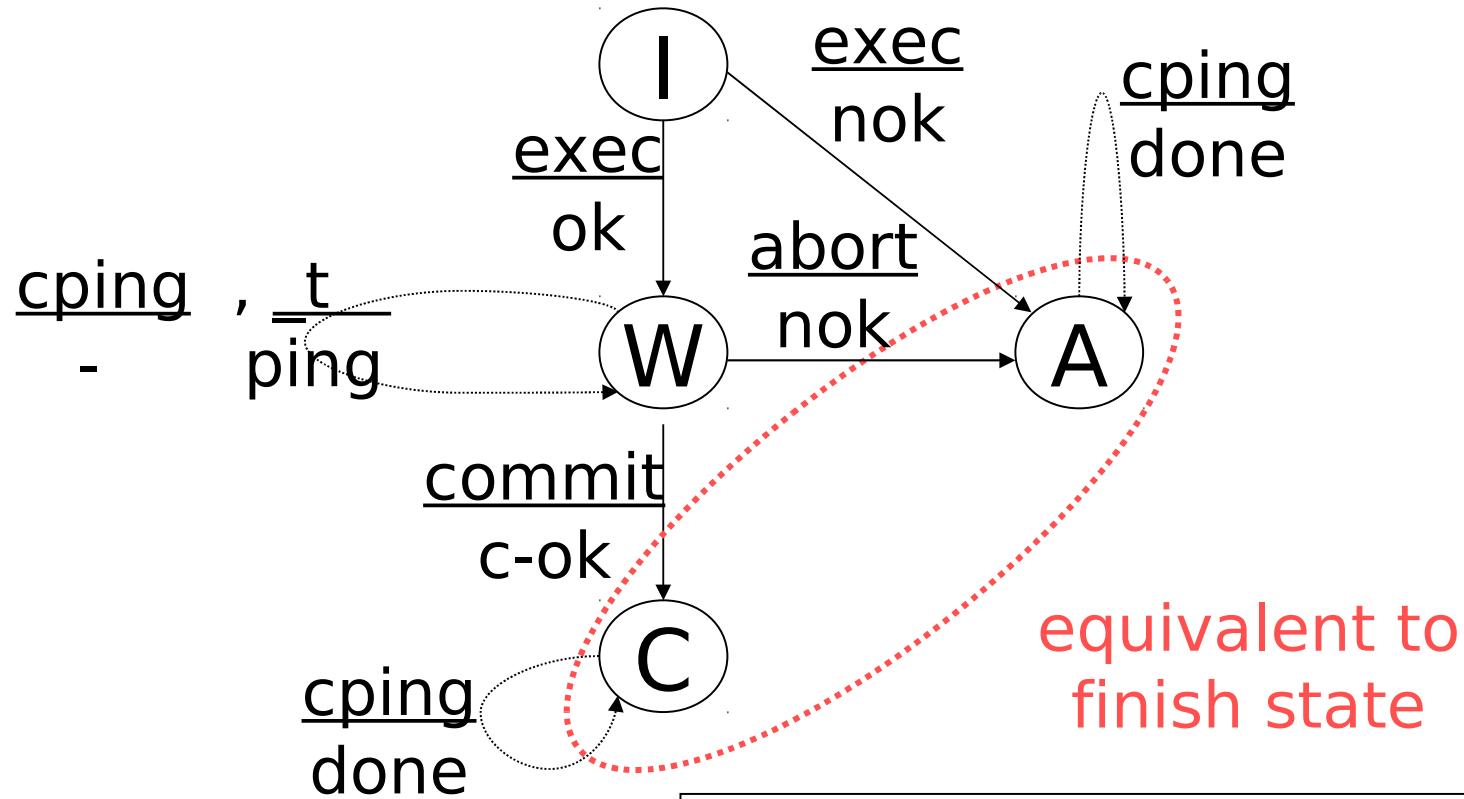


Participant



“done” message counts as either c-ok or n-ok for coordinator

Participant

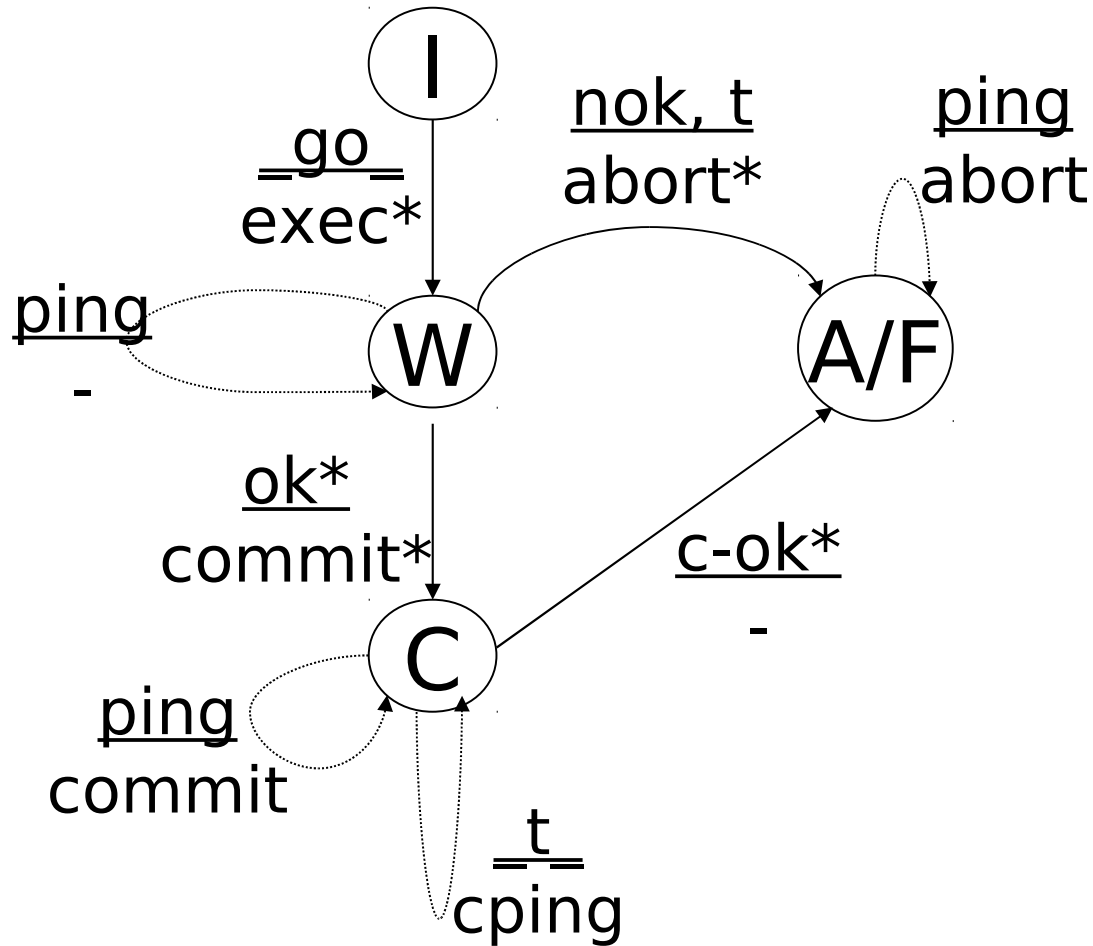


“done” message counts as either c-ok or n-ok for coordinator

Presumed abort protocol

- “F” and “A” states combined in coordinator
- Saves persistent space (forget quicker)
- Presumed commit is analogous

Presumed abort-coordinator (participant unchanged)



Remember:

all state transitions must be logged

Example: tracking who has sent “OK” msgs

Log at coord:

...	T1 start part={ <i>a,b</i> }	...	T1 OK from <i>a</i> RCV	...
-----	------------------------------------	-----	----------------------------------	-----

- After failure, we know still waiting for OK from node *b*
- Alternative: do not log receipts of “OK”s abort T1

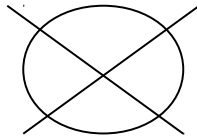
Example: logging receipt of C-OK messages

- If logged, can recover state
- If not logged:
 - resend commit *
 - participants reply “done” if duplicate

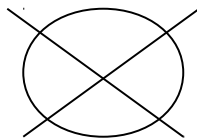
2PC is blocking

Sample scenario:

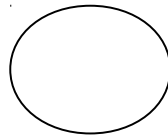
Coord



P1

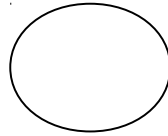


P2



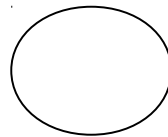
W

P3



W

P4



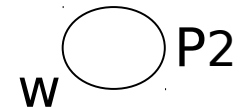
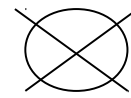
W

Case I:

$P_1 \rightarrow "W";$ coordinator sent commits

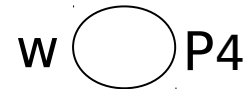
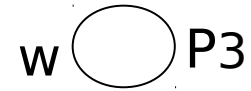
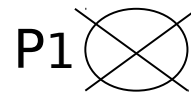
$P_1 \rightarrow "C"$

coord



Case II:

$P_1 \rightarrow \text{NOK}; P_1 \rightarrow A$

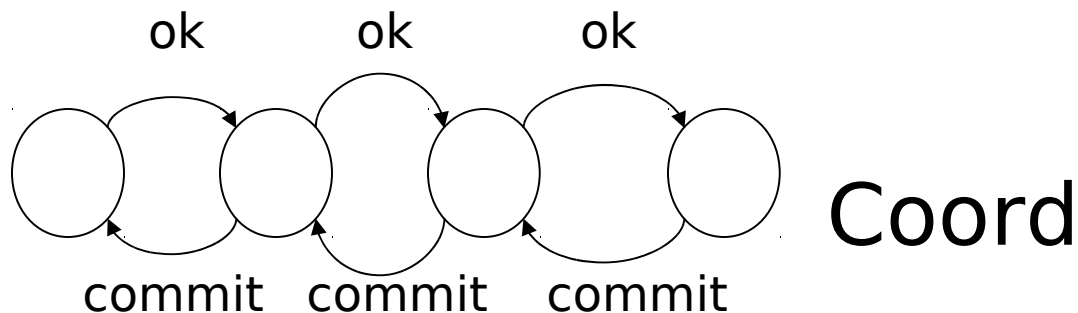


$\Rightarrow P_2, P_3, P_4$ (surviving participants)

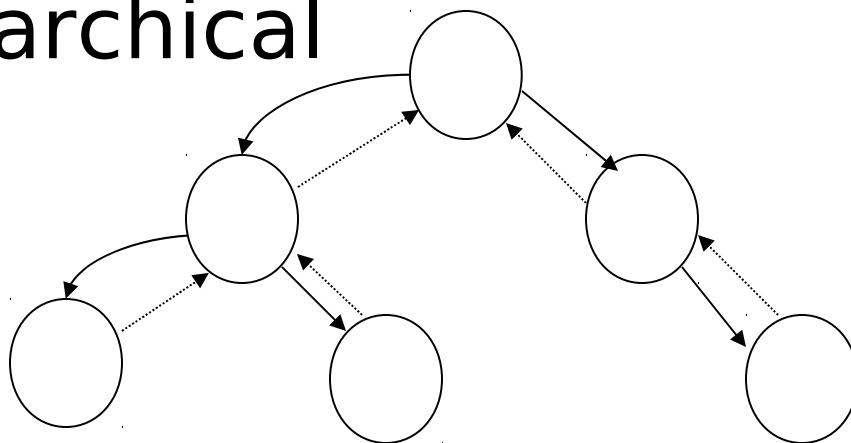
cannot safely abort or commit transaction

Variants of 2PC

- Linear

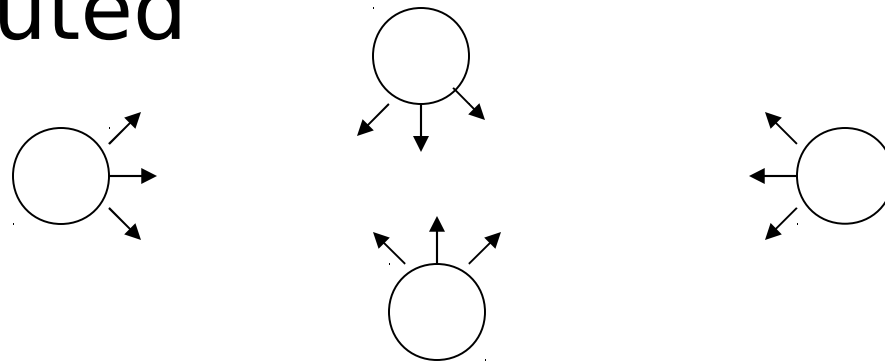


- Hierarchical



Variants of 2PC

- Distributed



- Nodes broadcast all messages
- Every node knows when to commit

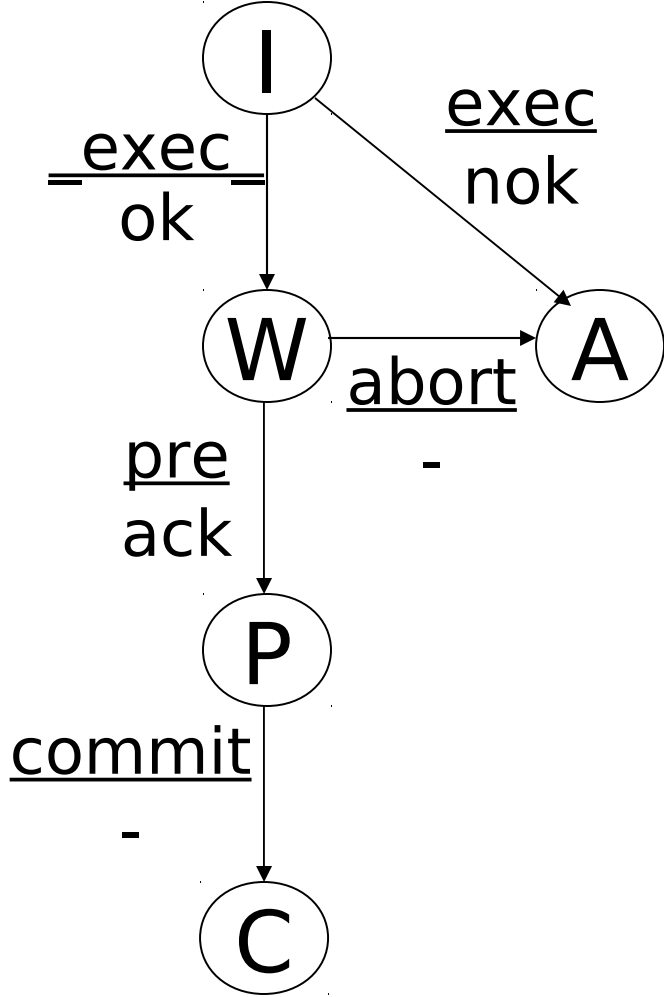
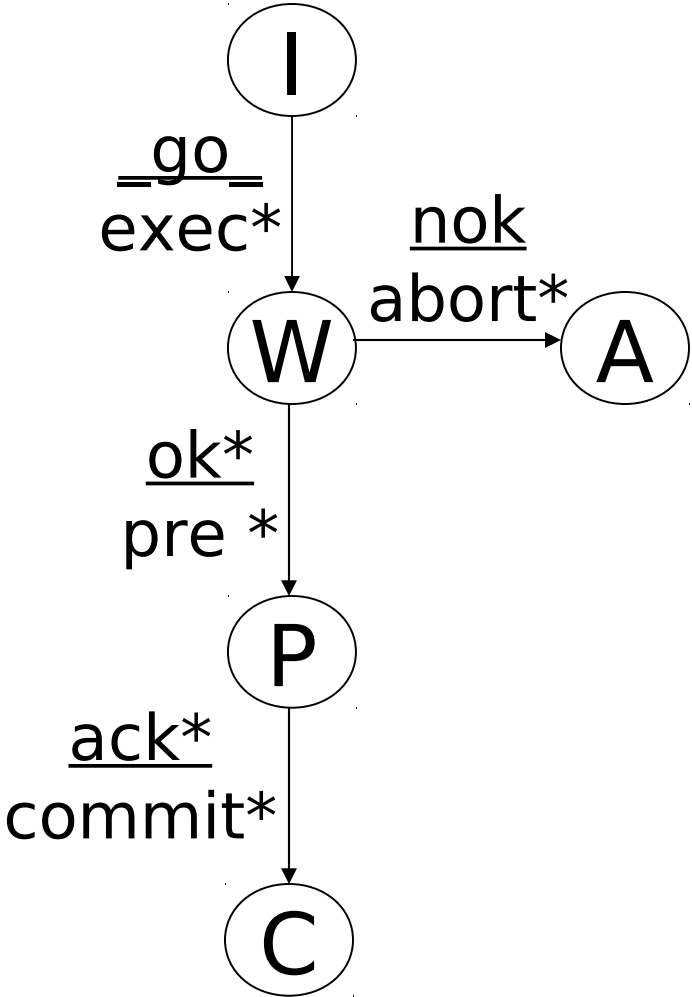
3PC = non-blocking commit

- Assume: failed node is down forever
- Key idea: before committing, coordinator tells participants everyone is ok

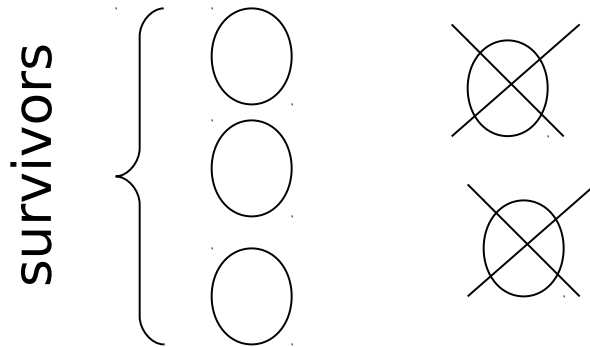
3PC

Coordinator

Participant

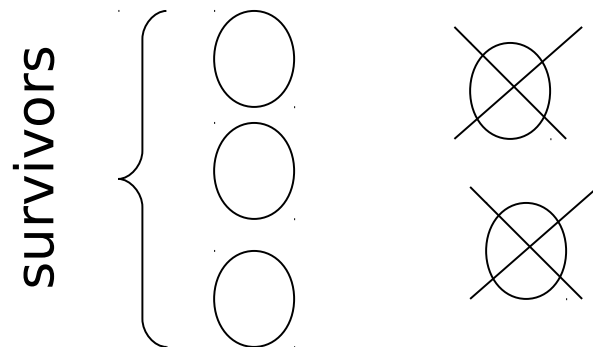


3PC recovery rules: termination protocol



- Survivors try to complete transaction, based on their current states
- Goal:
 - If dead nodes committed or aborted, then survivors should not contradict!
 - Else, survivors can do as they please...

- Let $\{S_1, S_2, \dots, S_n\}$ be survivor sites
- If one or more $S_i = \text{COMMIT} \Rightarrow \text{COMMIT } T$
- If one or more $S_i = \text{ABORT} \Rightarrow \text{ABORT } T$
- If one or more $S_i = \text{PREPARE} \Rightarrow$
 T could not have aborted $\Rightarrow \text{COMMIT } T$
- If no $S_i = \text{PREPARE}$ (or COMMIT) \Rightarrow
 T could not have committed $\Rightarrow \text{ABORT } T$



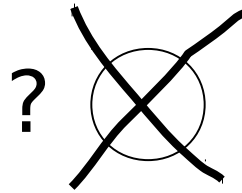
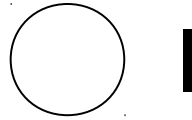
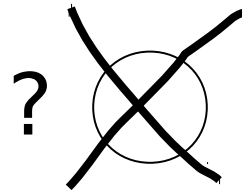
Example:

? ~~○~~ ○ P

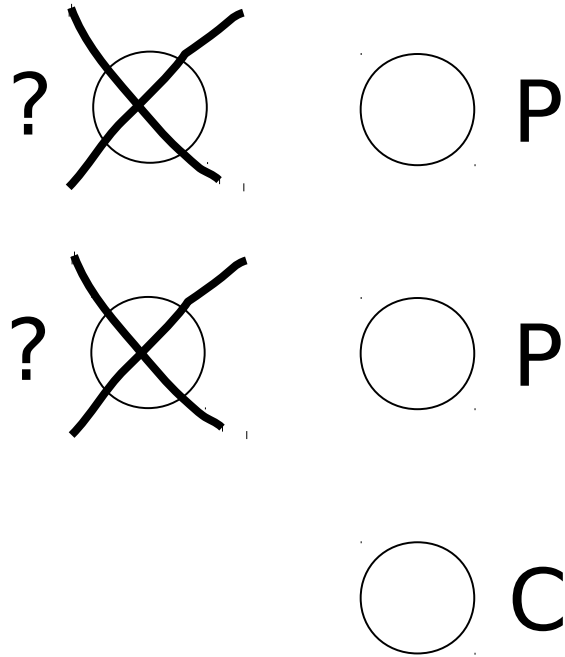
? ~~○~~ ○ W

○ W

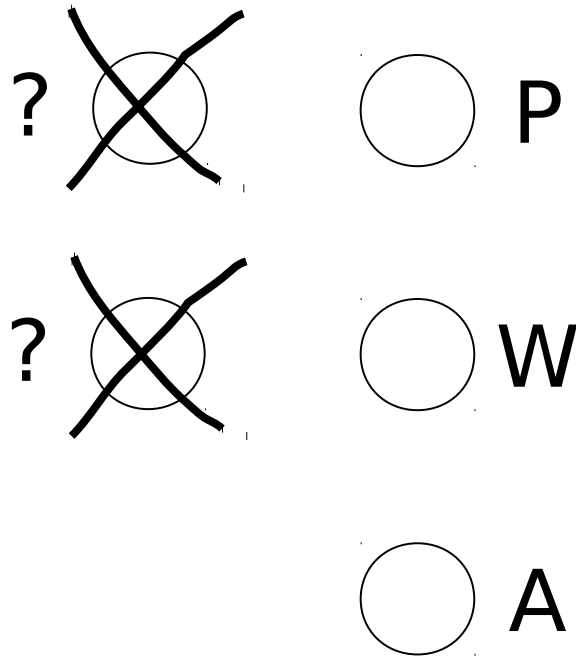
Example:



Example:



Example:



⇒ Once survivors make decision, they must

select new coordinator to continue

3PC

P ○

W ○

W ○

Decide to
commit

Time
1

P ○

P ○

P ○

Time
2

C ○

C ○

P ○

Time
3

C ○

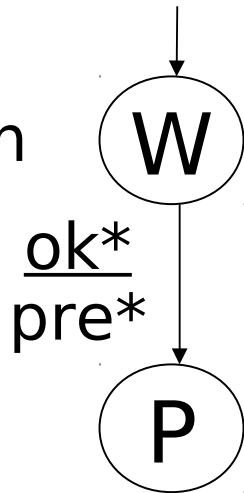
C ○

C ○

Time
4

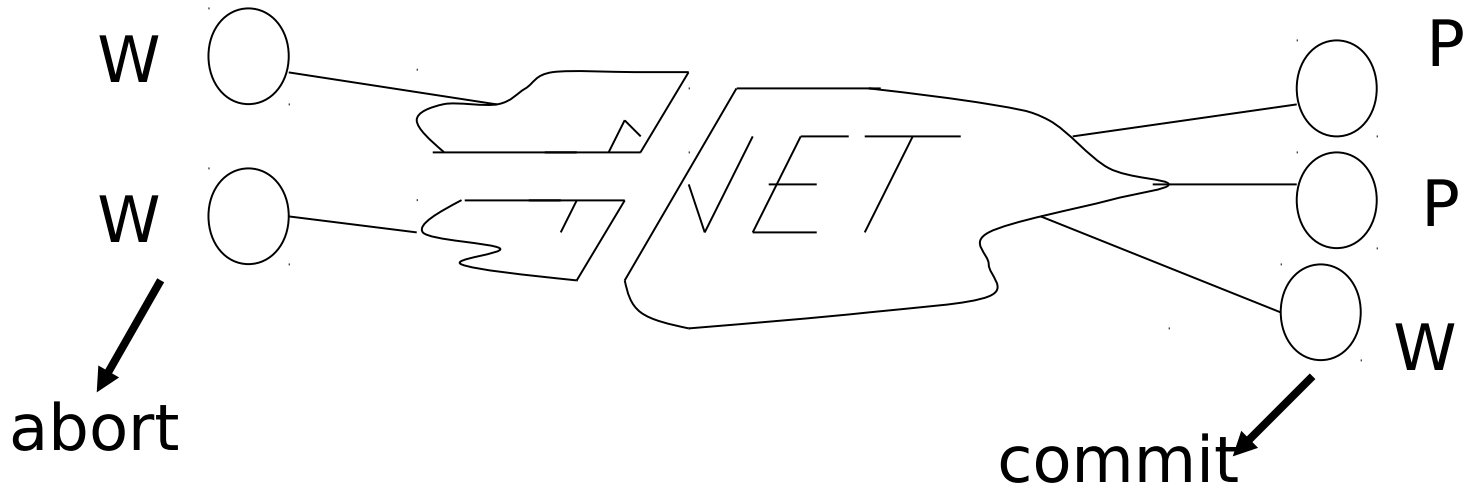
Note: when survivors continue 3PC,
failed nodes do not count

E.g.,
from



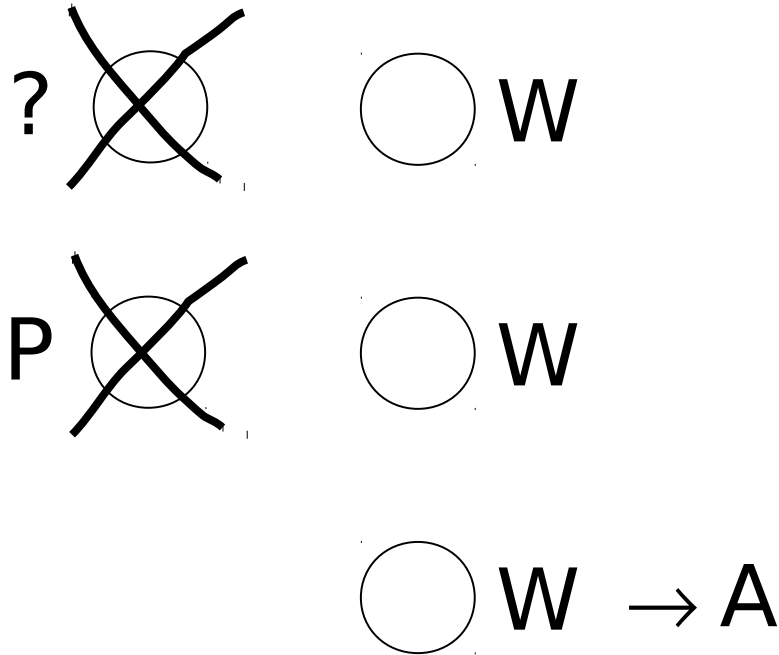
“OK*” ≡ when OK’s received
non-failed nodes

Note: 3PC unsafe with partitions!



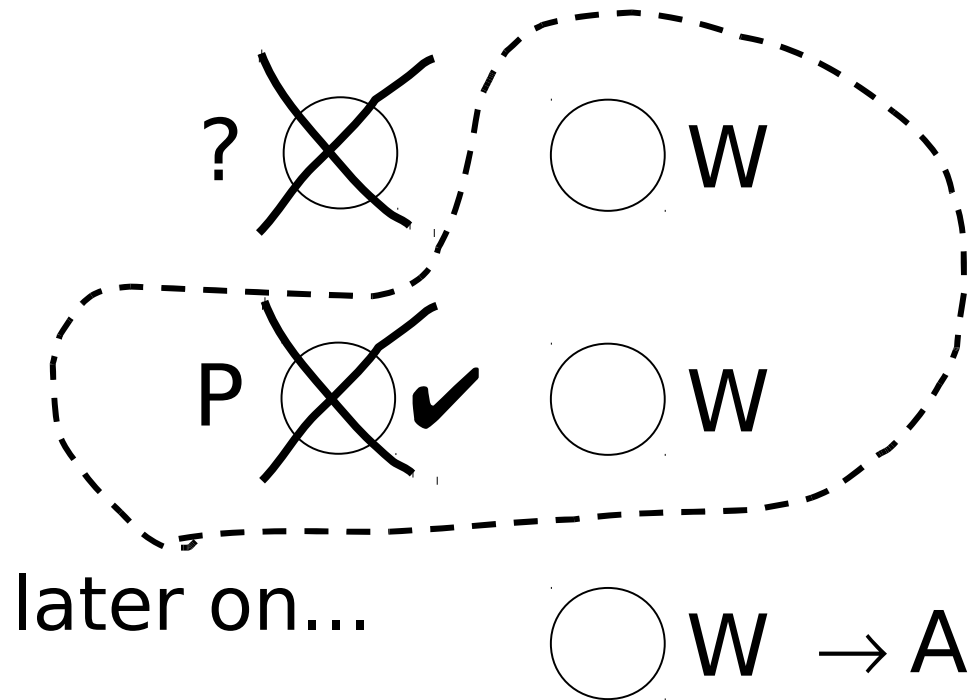
Node recovery:

- After node N recovers from failure:
 - do not participate in termination protocol (why?)



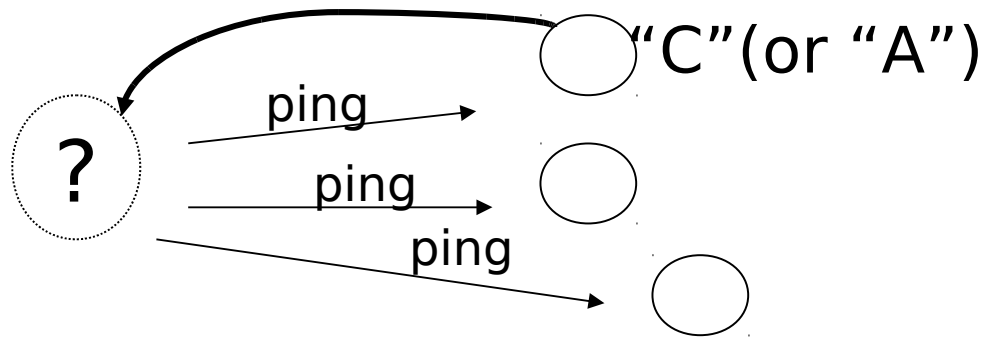
Node recovery:

- After node N recovers from failure:
 - do not participate in termination protocol (why?)



Node recovery:

- After node N recovers from failure:
 - do not participate in termination protocol (why?)
 - wait until it hears commit or abort decision from operational node



- Waiting for commit/abort decision from other node is ok, unless all fail:



Two options for all-failed problem:

- (A) Wait for all to recover
- (B) Majority commit

Option A

- Recovering node waits for either:
 - (1) commit/abort outcome for T from other node
 - (2) all nodes that participated in T are up and recovering:
 - ⇒ then 3PC can continue
 - (no danger that a failed node could have aborted or committed)

Example(1): Coord ~~⊗~~ ~~⊗~~ - ~~⊗~~ ? ~~⊗~~ V

P₁ P₃ → W ~~⊗~~ ?

p₄ → W ?

- Nodes P₂, P₃, P₄ enter “W” state and fail
- When they recover, coord. and P₁ are down
- Each node has 1 vote, V=5, Maj=3

Example(2): Coord ~~⊗~~ ~~⊗~~ P₃ → "P"
 P₁ ~~⊗~~ ~~⊗~~ P₄ → "W"
 P₂ ~~⊗~~

- Each node has 1 vote; $V=5$, $Maj=3$
- Nodes fail after entering states shown;
 P₃, P₄ recover

Example(2): Coord ~~⊗~~ ~~⊗~~ P₃ → "P"
 P₁ ~~⊗~~ ~~⊗~~ P₄ → "W"
 P₂ ~~⊗~~

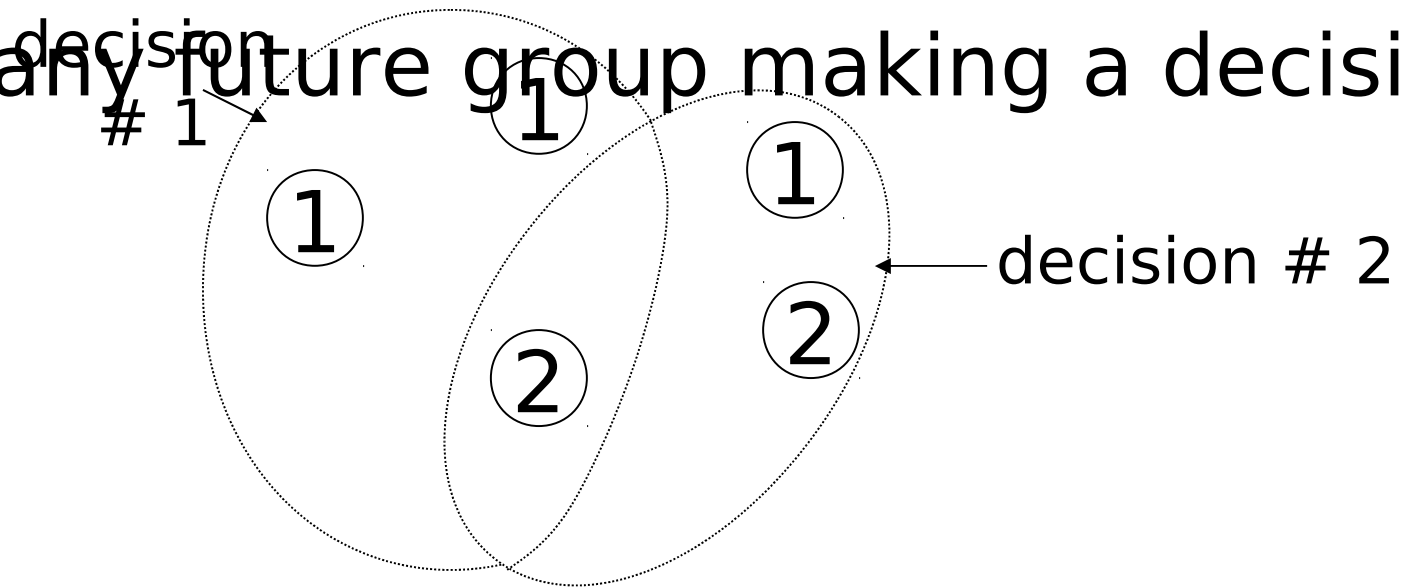
- Each node has 1 vote; $V=5$, $Maj=3$
- Nodes fail after entering states shown;
 P₃, P₄ recover
- Termination rule says we can try to commit, but P₃, P₄ do not have enough votes, so they do nothing!
- P₃, P₄ doing nothing is good because later on, coord. P₁, P₂ could abort T

Summary:

Majority rule ensures that any decision

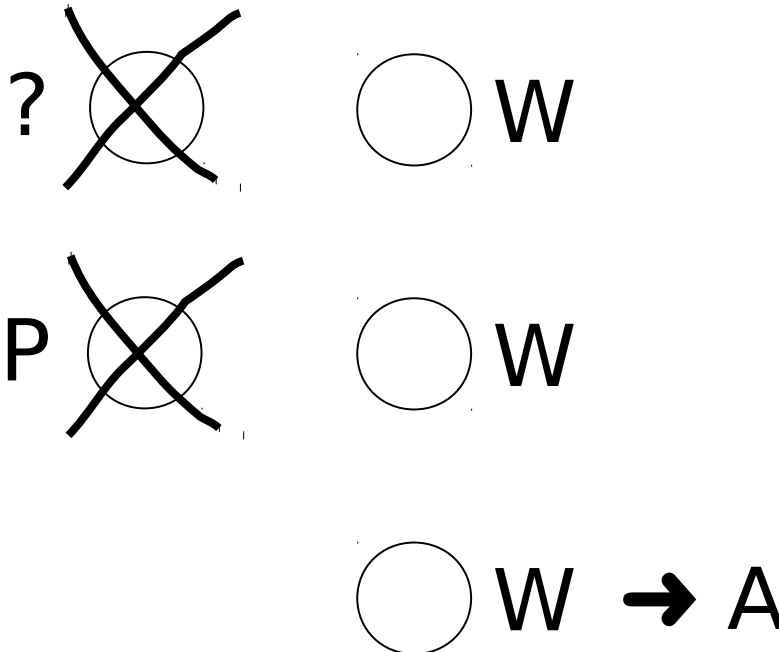
(e.g., Preparing, committing) will be known

to any future group making a decision



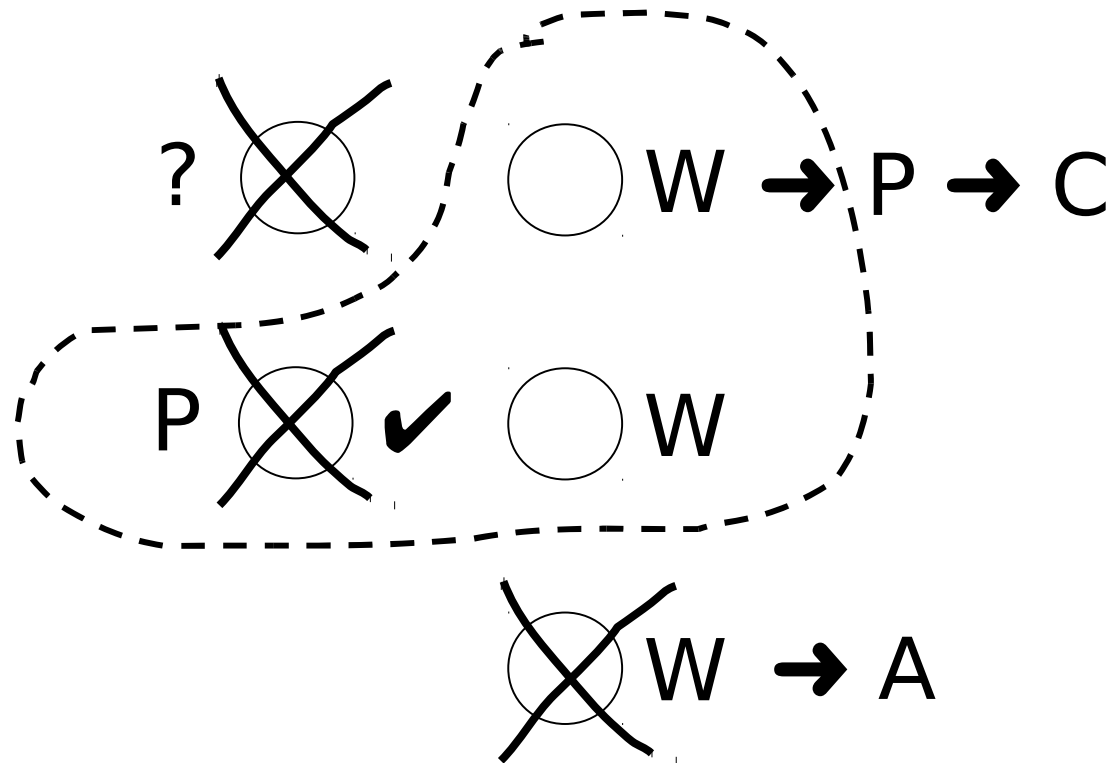
Important Detail for Majority 3PC

- Example:



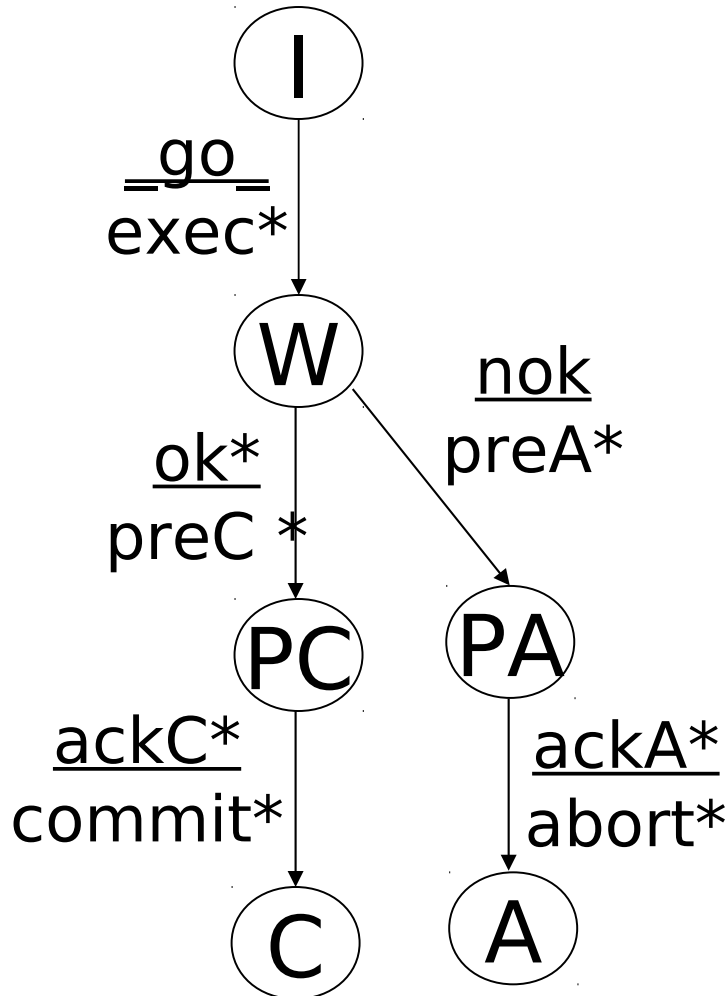
Important Detail for Majority 3PC

- Example:

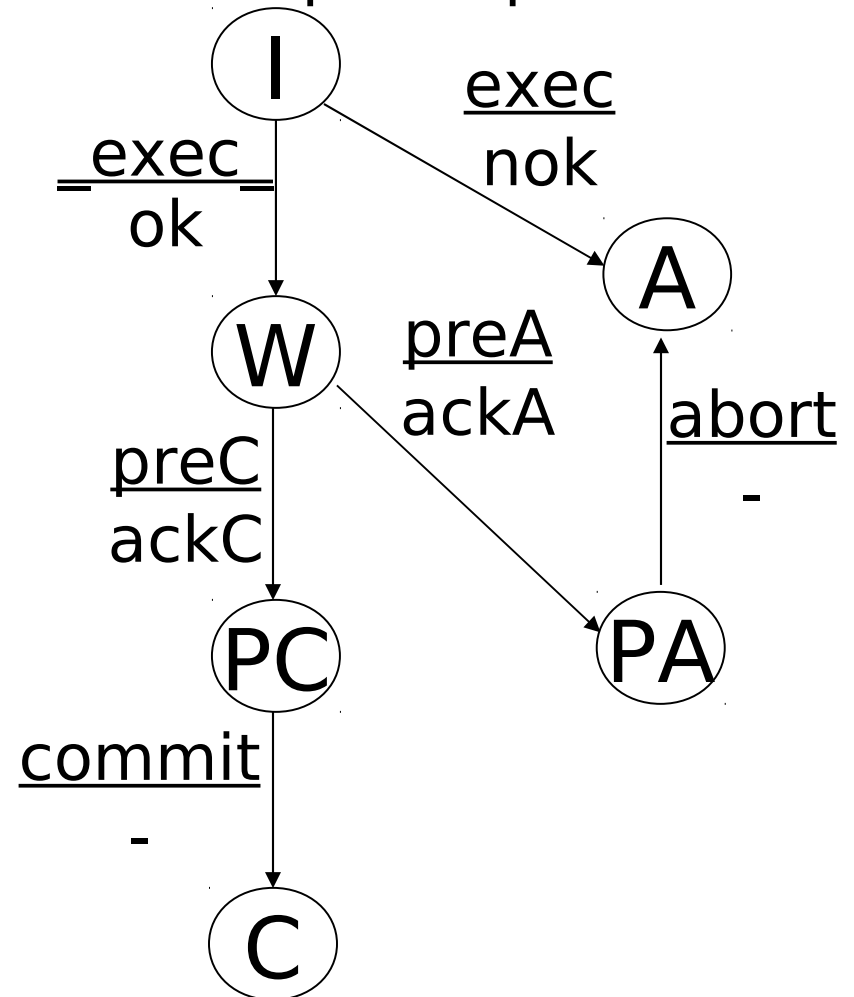


Need "Prepare To Abort" State

coordinator



participant



Example Revisited

? ~~○~~

○ W

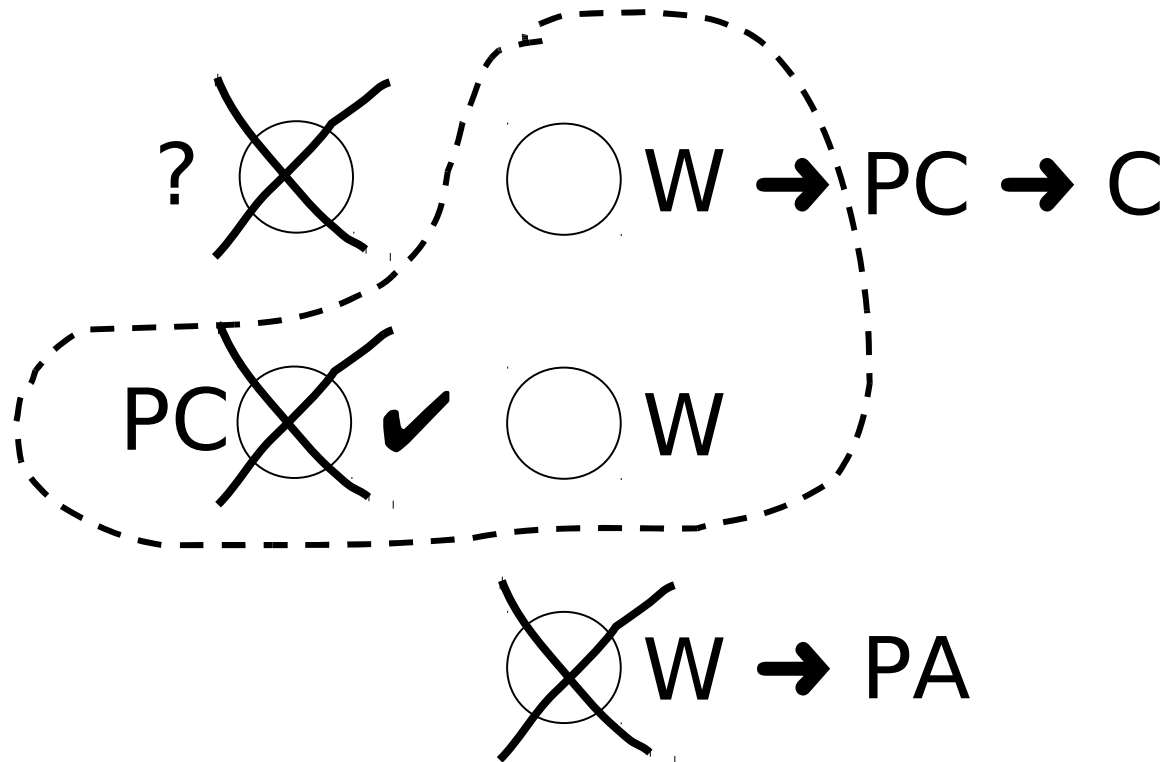
PC ~~○~~

○ W

○ W → PA

Example Revisited

OK to commit since
transaction could not have aborted



Example Revisited -II

? ~~○~~

○ W

PC ~~○~~

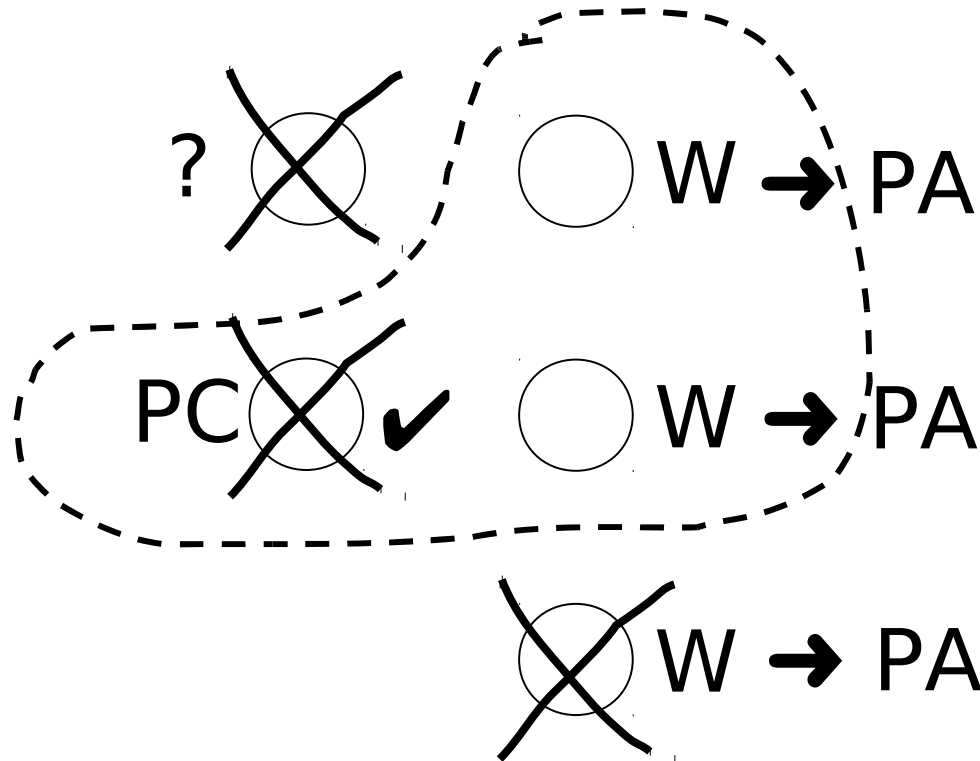
○ W → PA

○ W → PA

Example Revisited -II

No decision:

Transaction could have aborted or
could have committed... Block!



3PC with Majority Voting

- If survivors have majority and all states $W \Rightarrow$ try to abort
- If survivors have majority and states in $\{W, PC, C\} \Rightarrow$ try to commit
- If survivors have majority and states in $\{W, PA, A\} \Rightarrow$ try to abort
- Otherwise block

Comparison

Option A: only nodes that have not failed participate in 3PC

- Any size group can terminate
(even one node)
- If all nodes fail, must wait for all to recover

Comparison

Option B: Majority voting

- A group of failed+recovering nodes can terminate transaction (with majority of votes)
- Need majority for every commit
 - blocking protocol!

Reminder

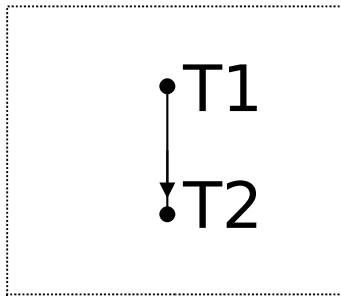
- When node recovers, it uses its log in a normal fashion to determine status of transactions:
 - if commit found in log \Rightarrow redo if necessary
 - if abort found (or no “W” record) \Rightarrow rollback if necessary

Reminder - Continued

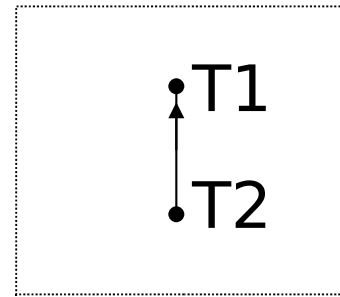
- if in “W” state (or “P” state):
 - reclaim locks held by T before crash
 - try to terminate T (with other nodes)
- after locks claimed for “in doubt” transactions, start normal processing

Final note

- If nodes use 2P locking, global deadlocks possible

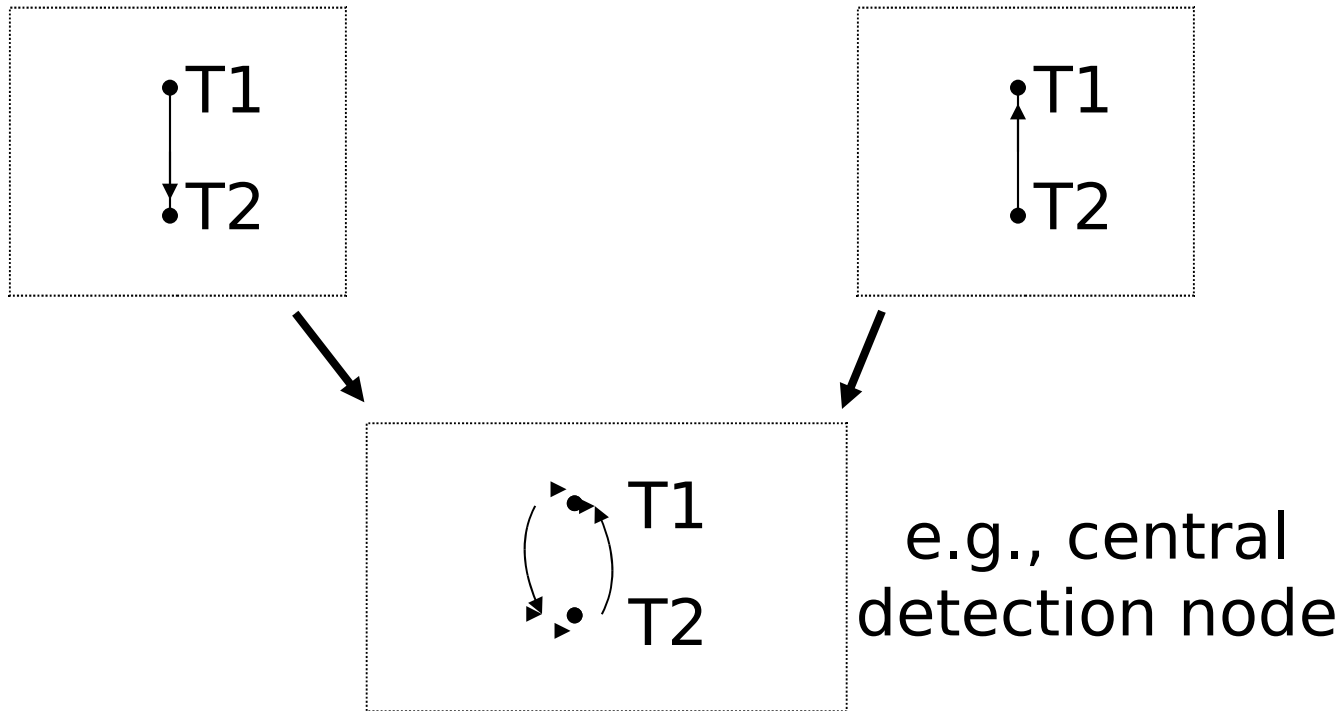


Local WFG:
no cycles

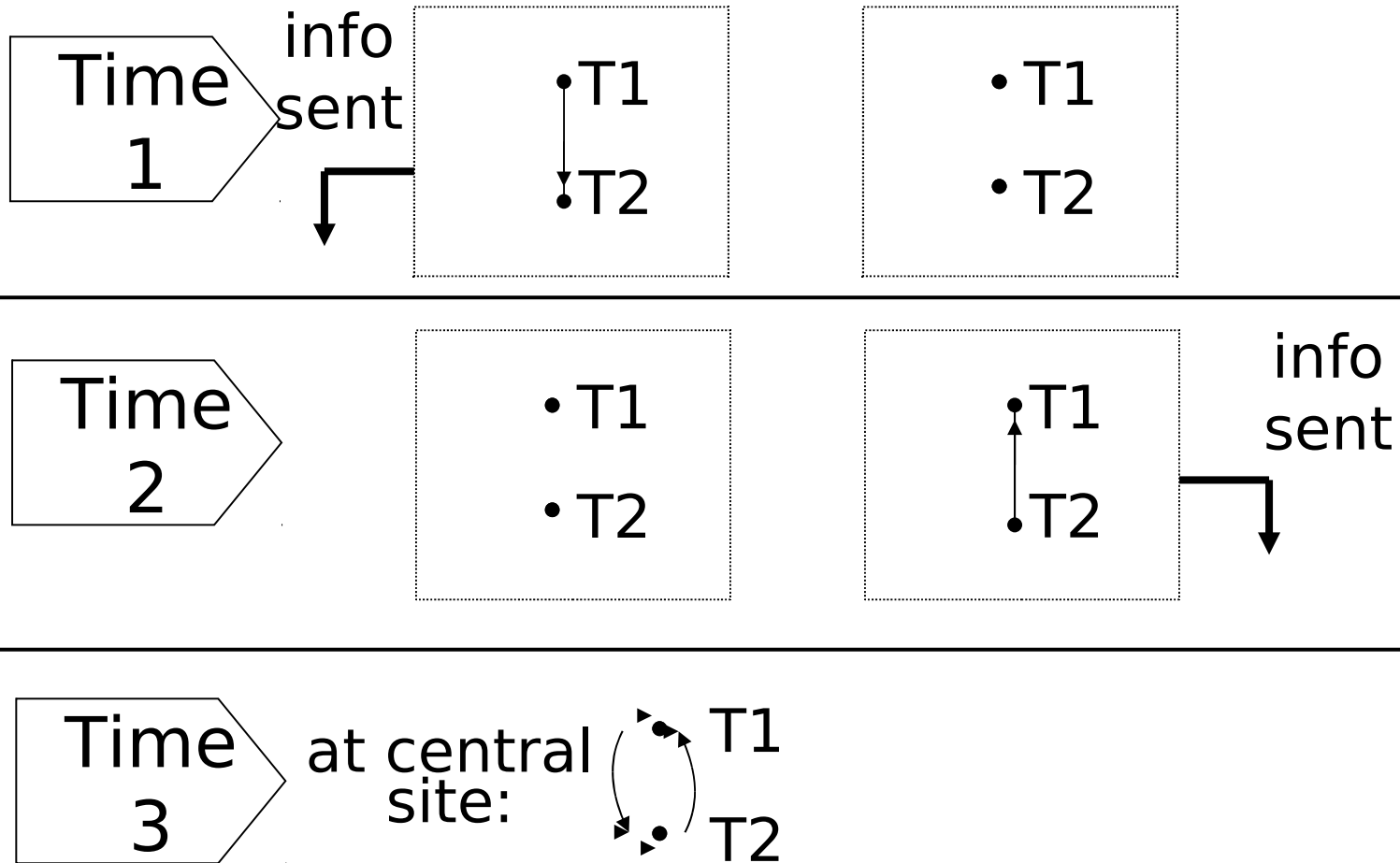


Local WFG:
no cycles!

- Need to “combine” WFGs to discover global deadlock



Problem: False deadlocks



- Many deadlock solutions
 - Distributed vs. centralized
 - Detection vs. prevention
 - timeouts
 - wait-die
 - wound-wait
- Covered in CS245