# CS 3351: Paxos Made Simple?

Amy Babay

University of Pittsburgh
School of Computing and Information
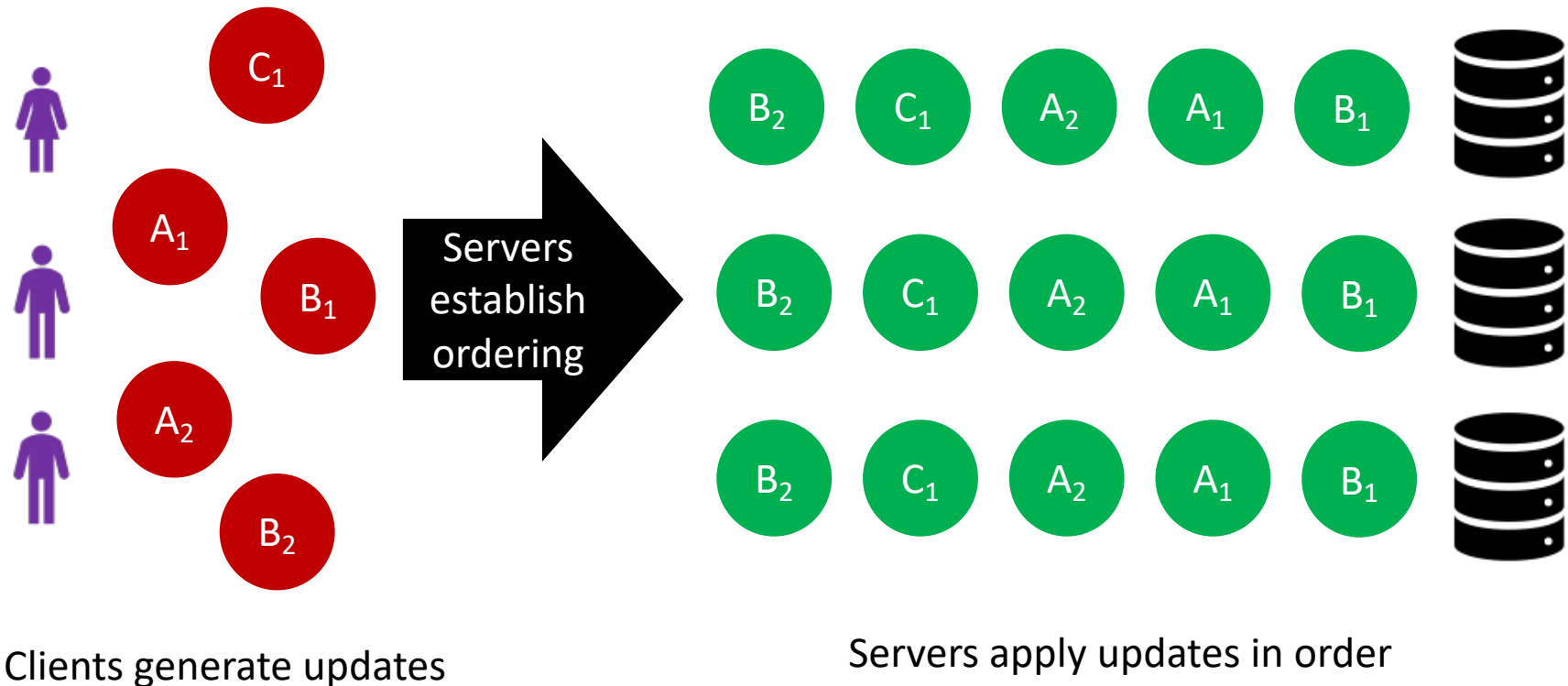
# Background & Motivation

- **State Machine Replication (SMR)**: technique for implementing strongly consistent fault-tolerant services
  - Servers start in the same state
  - Servers apply deterministic updates in the same order
  - => Servers progress through exactly the same sequence of states

# Background & Motivation

- **State Machine Replication (SMR)**: technique for implementing strongly consistent fault-tolerant services



Clients generate updates

Servers apply updates in order

# Background & Motivation

- As discussed last class, SMR can be implemented as a sequence of **consensus** instances
- **Paxos** is a consensus protocol that is often used to implement SMR
  - Published in technical report in 1989, Journal paper in 1998 ("The Part Time Parliament")
    - Described via analogy to hypothetical Greek parliament
  - Work on clarifying, implementing, optimizing, or replacing Paxos continues in the literature today
  - Foundation for some of the intrusion-tolerant replication protocols we'll see later
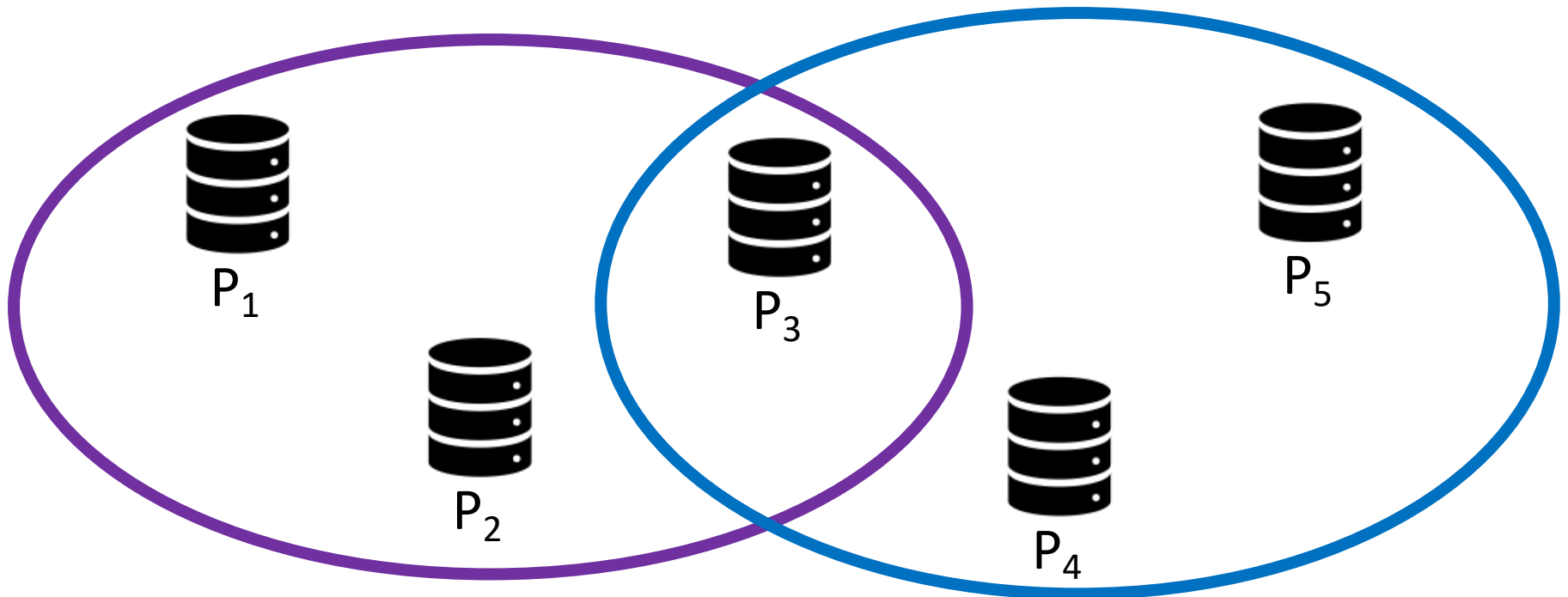
# System Model

- We have a groups of processes (aka. servers, replicas)
- Processes communicate by sending messages
- Processes can crash *and restart*
  - Processes have access to stable storage to record information
- Messages can take arbitrarily long to be delivered, can be duplicated, and can be lost
- Processes execute the protocol faithfully, and messages are not corrupted (non-Byzantine)

# Consensus Requirements (Safety)

- **Validity**: Only a value that has been proposed may be chosen

- **Agreement (1)**: Only a single value is chosen

- **Agreement (2)**: A process never learns that a value has been chosen unless it actually has been

- What about liveness (progress/termination)?

# Paxos Consensus Protocol ("Single-Decree Synod")

- Key concept: Any two majorities must intersect in at least one process

# Paxos Consensus Protocol ("Single-Decree Synod")

- Key concept: Any two majorities must intersect in at least one process

- So, to guarantee agreement we:

    1. Only allow a value to be *chosen* if it is accepted by a majority of (acceptor) processes

    2. Require a (proposer) process to communicate with a majority of (acceptor) processes before proposing a value to find out what they've previously accepted

# Paxos Consensus Protocol

- ## Subtle points:
  - But what if multiple values have been accepted?
  - Can a process accept a new value after giving its response?

- ## Solutions:
  - Sequence numbers on proposer's attempts to pass a proposal (often called "view numbers")
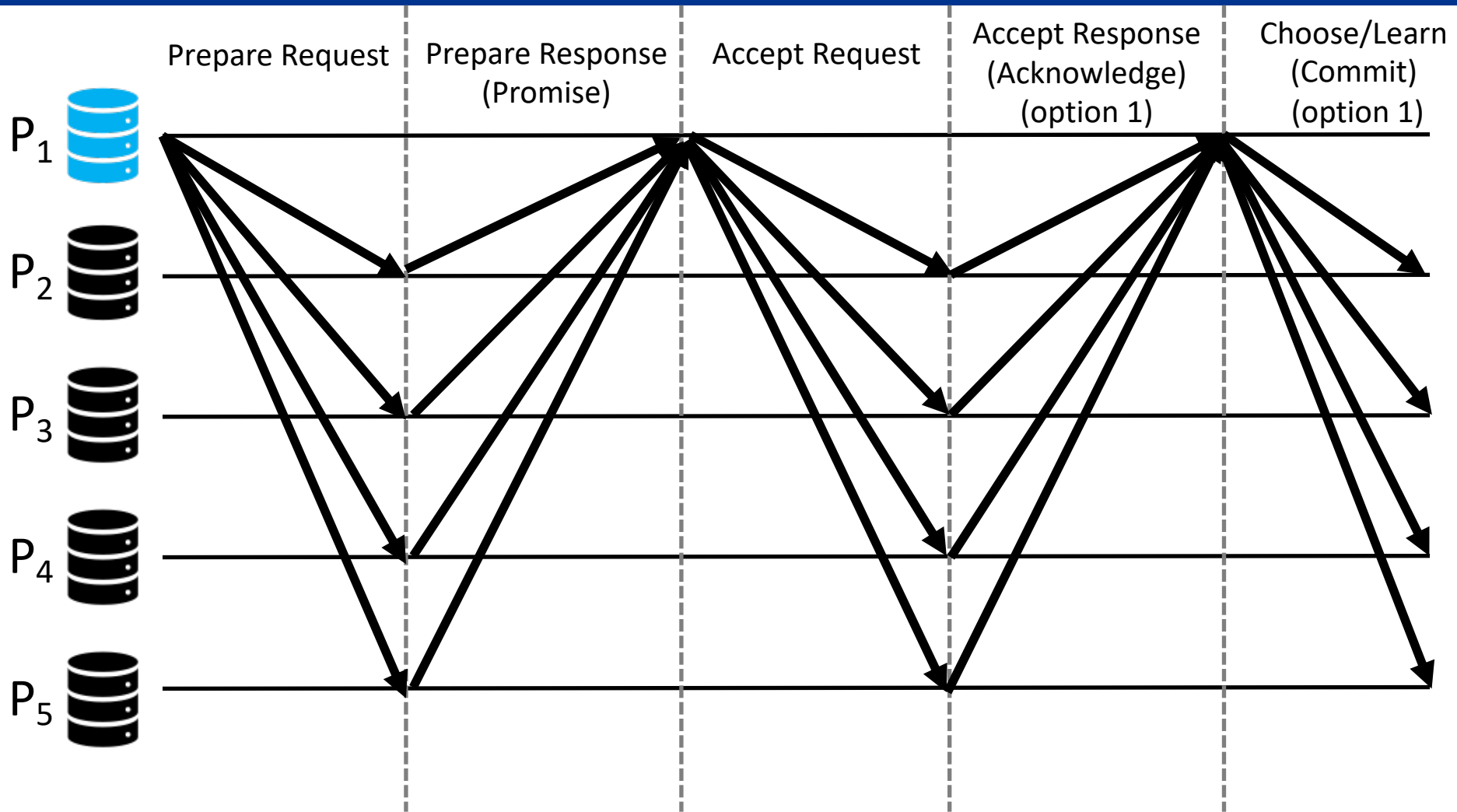  - Promises to not accept anything from lower views after responding

# Proposer Actions

- **Send** prepare_request(n)
- **Wait** for majority prepare_response(n,v,m)
  - v : value of last accepted proposal (may be null if none)
  - m: sequence number of last accepted proposal (m < n)
- *If* received some prepare_response(n,v,m) with non-null v, **send** accept_request(n,x) where x is the value associated with highest received m; *else* **send** accept_request(n,x), where x can be any value
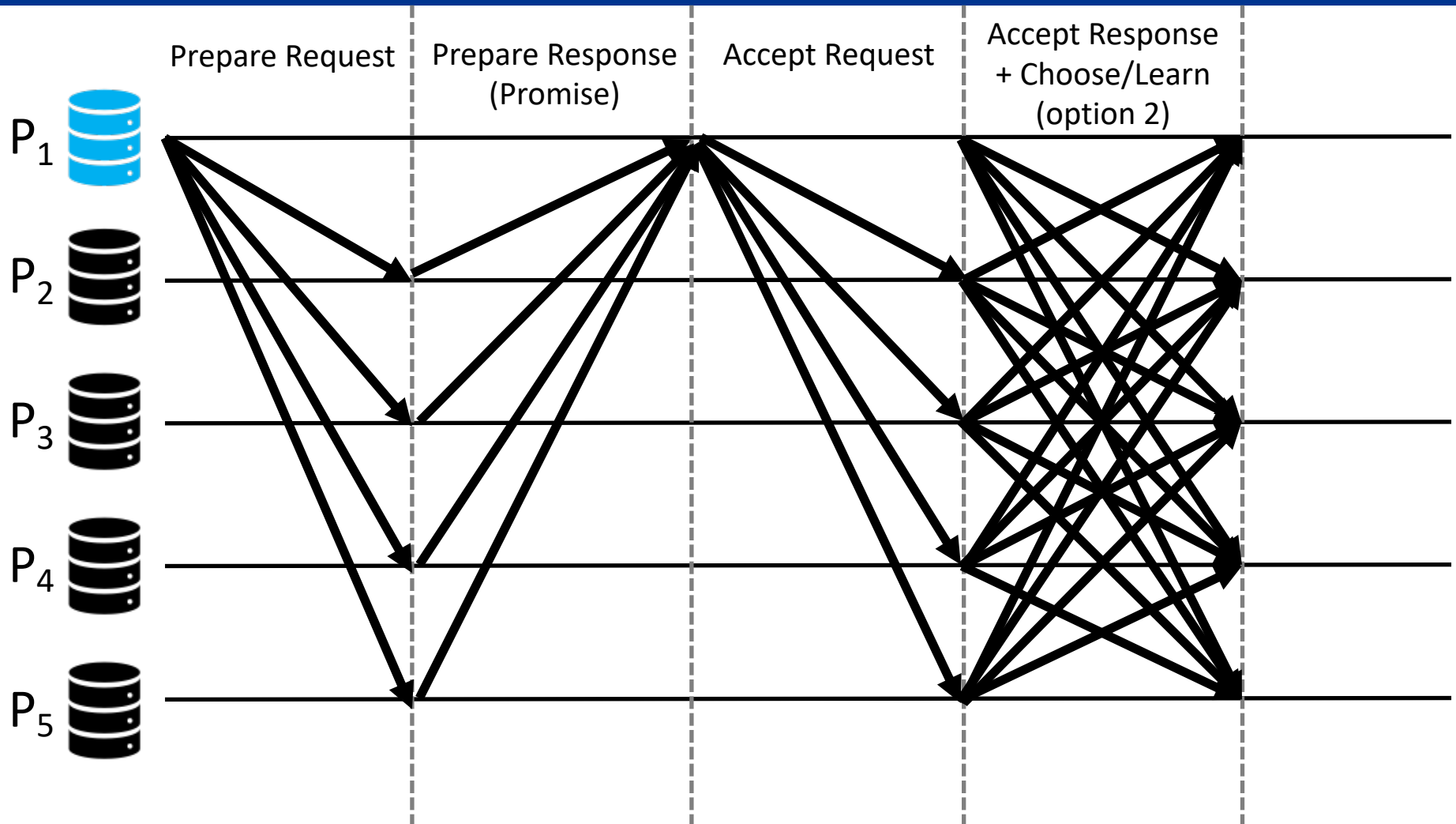
# Acceptor Actions

- **Upon receiving** prepare_request(n)
  - *If* previously responded to prepare_request(m) s.t. m > n, **do nothing** (or "inform proposer" as performance optimization)
  - *Else if* previously accepted prepare_request(m) s.t. m < n, **send** prepare_response(n,v,m)
  - *Else* **send** prepare_response(n,null,0)
- **Upon receiving** accept_request(n,x)
  - *If* previously responded to prepare_request(m) s.t. m > n, do nothing (or "inform proposer" as performance optimization)
  - *Else* **accept** value x (and inform learners)

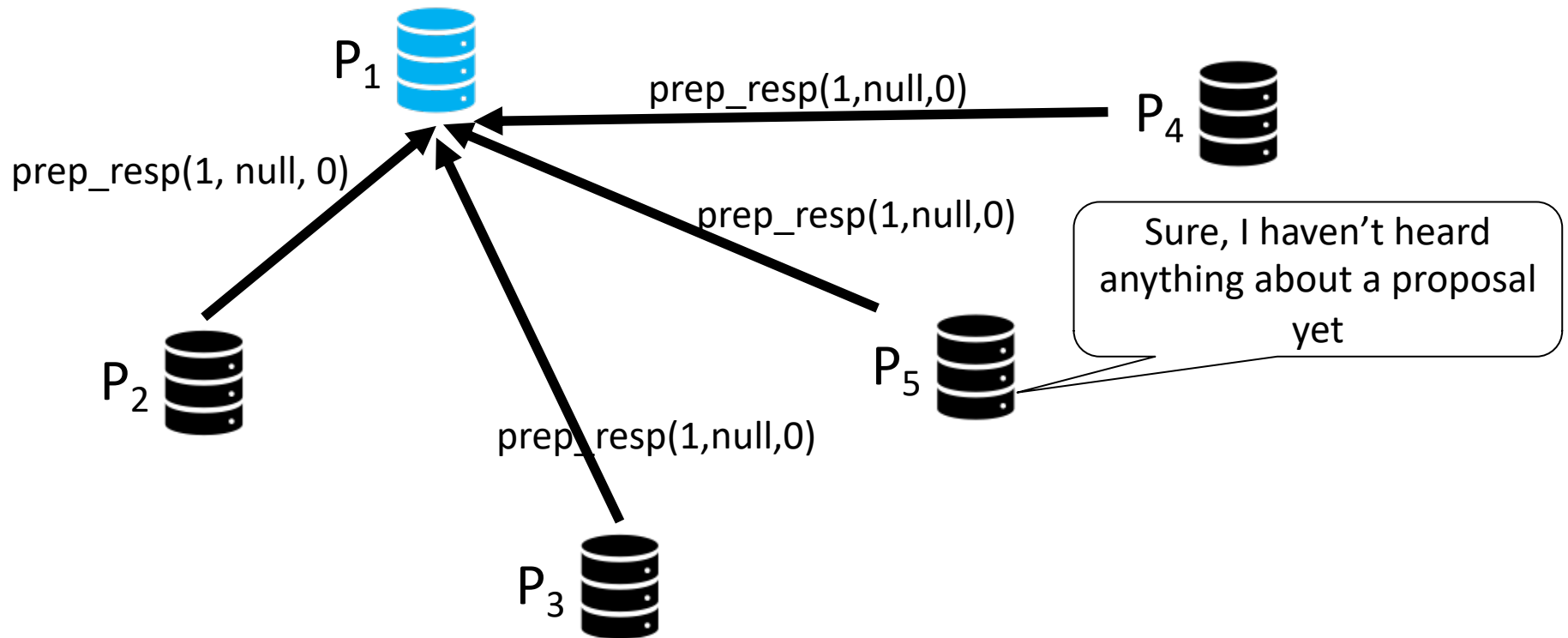# Paxos Consensus Protocol ("Single-Decree Synod")

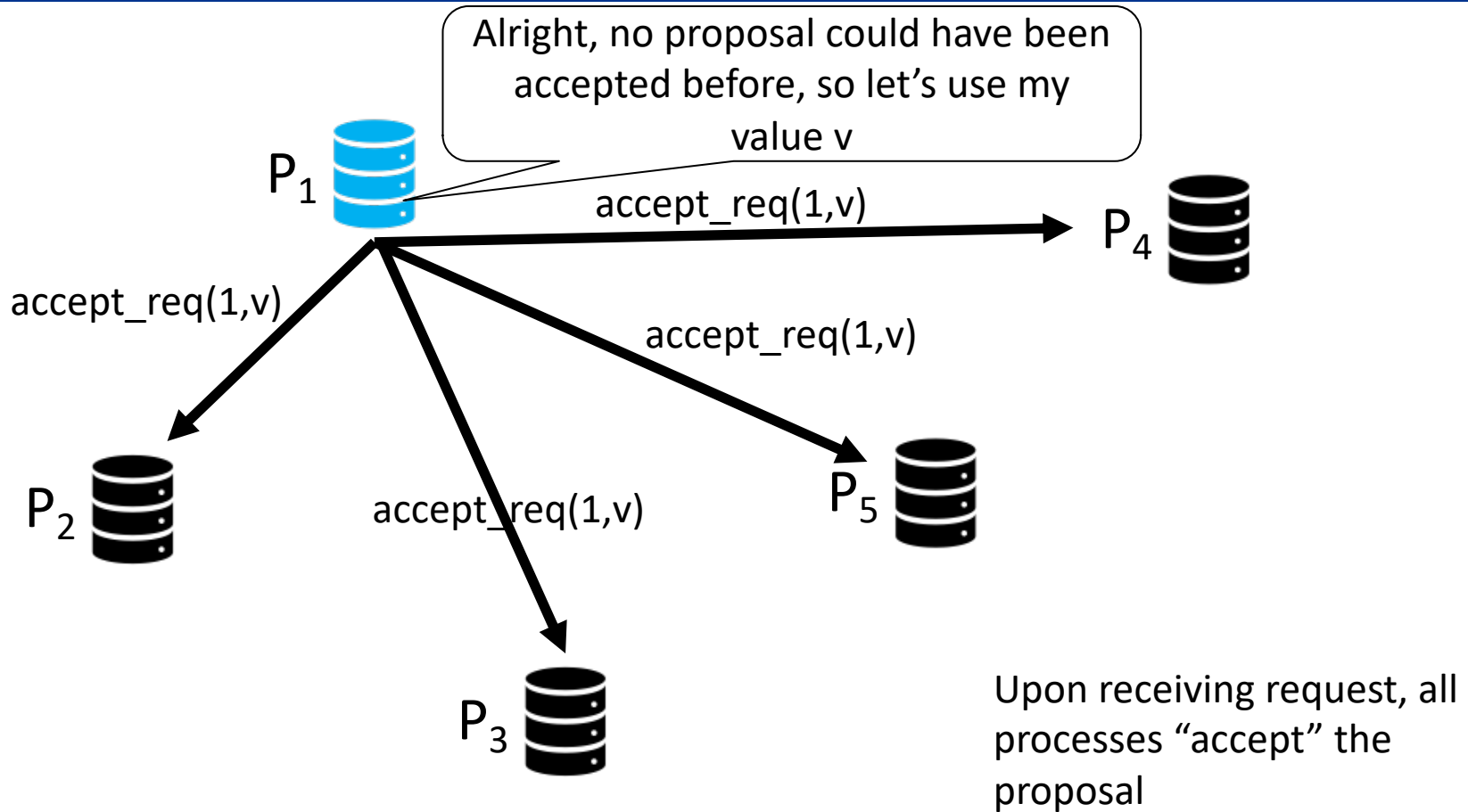# Paxos Consensus Protocol ("Single-Decree Synod")
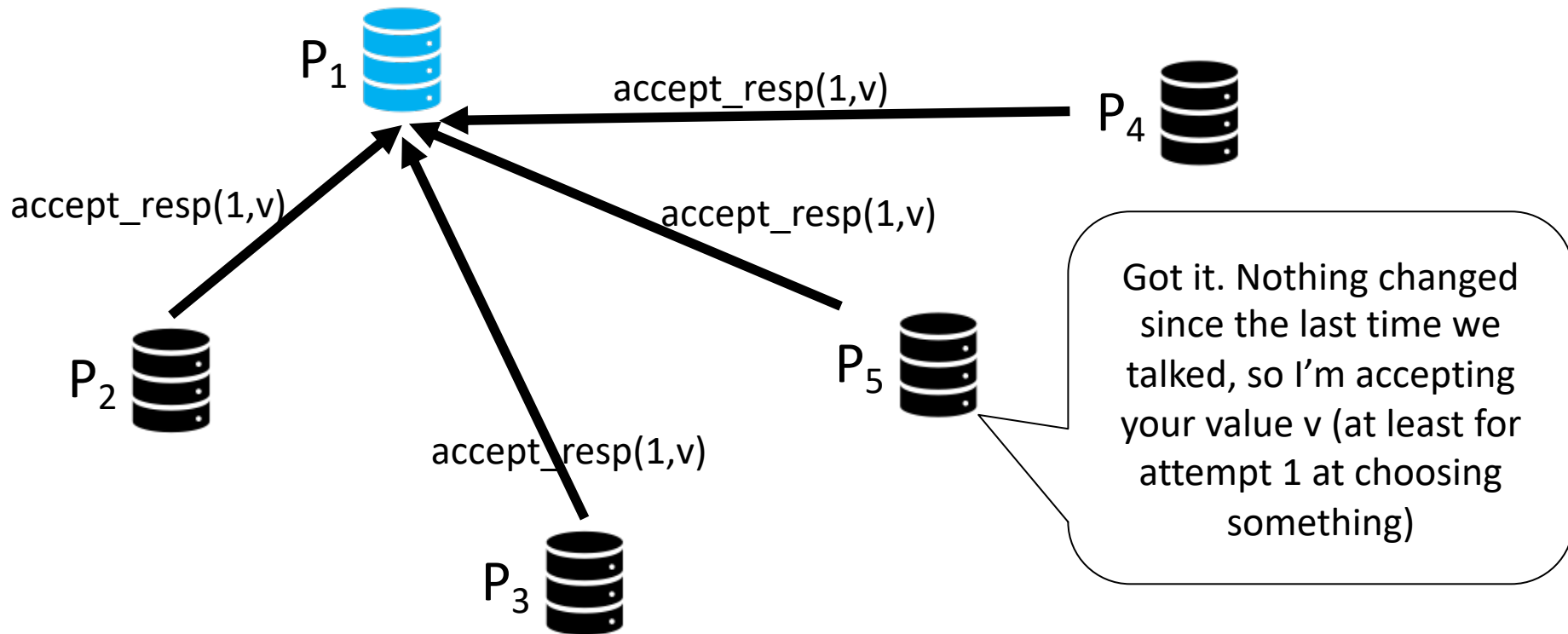
# Normal Case: Prepare Request

# Normal Case: Prepare Response

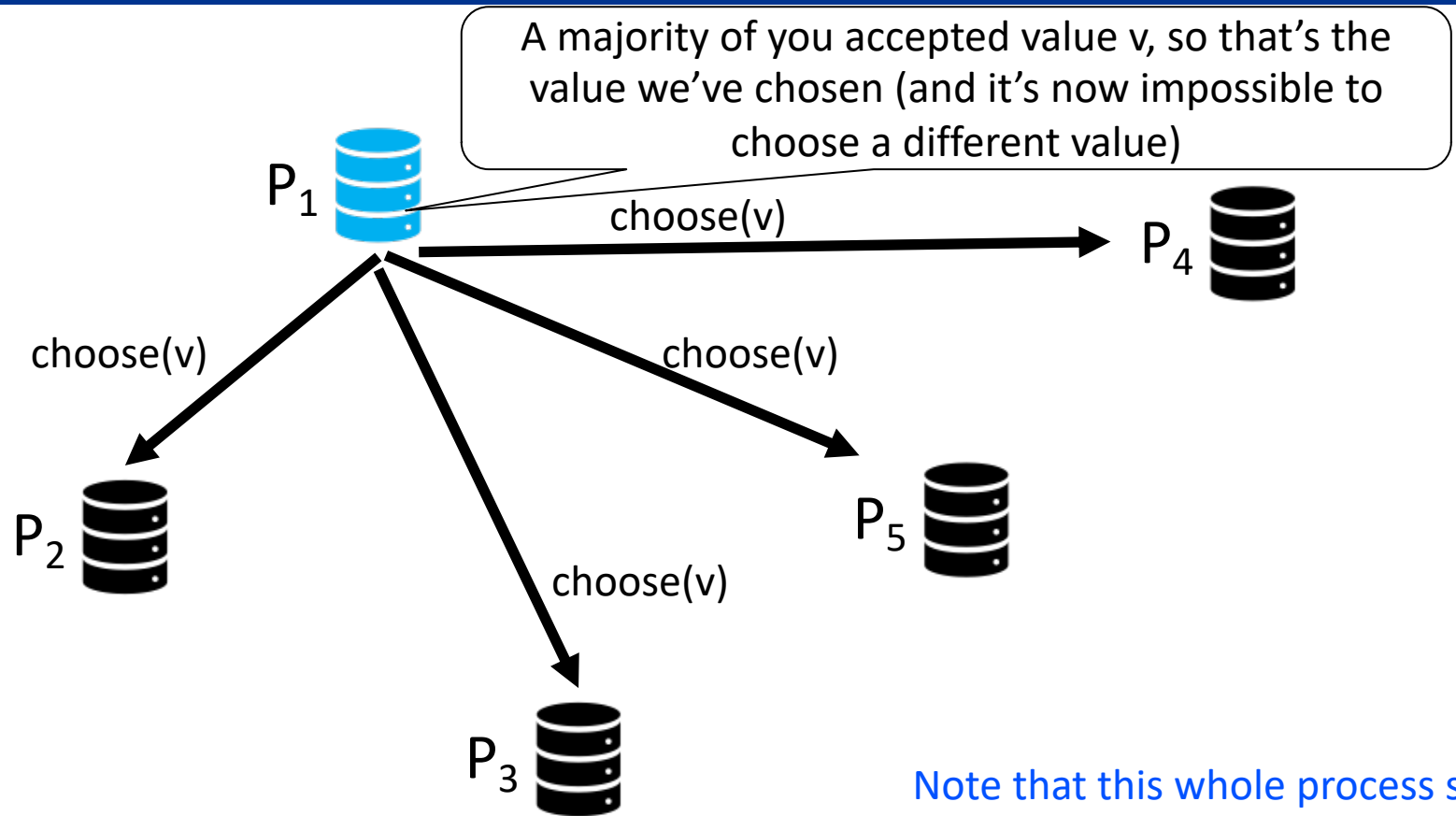# Normal Case: Accept Request
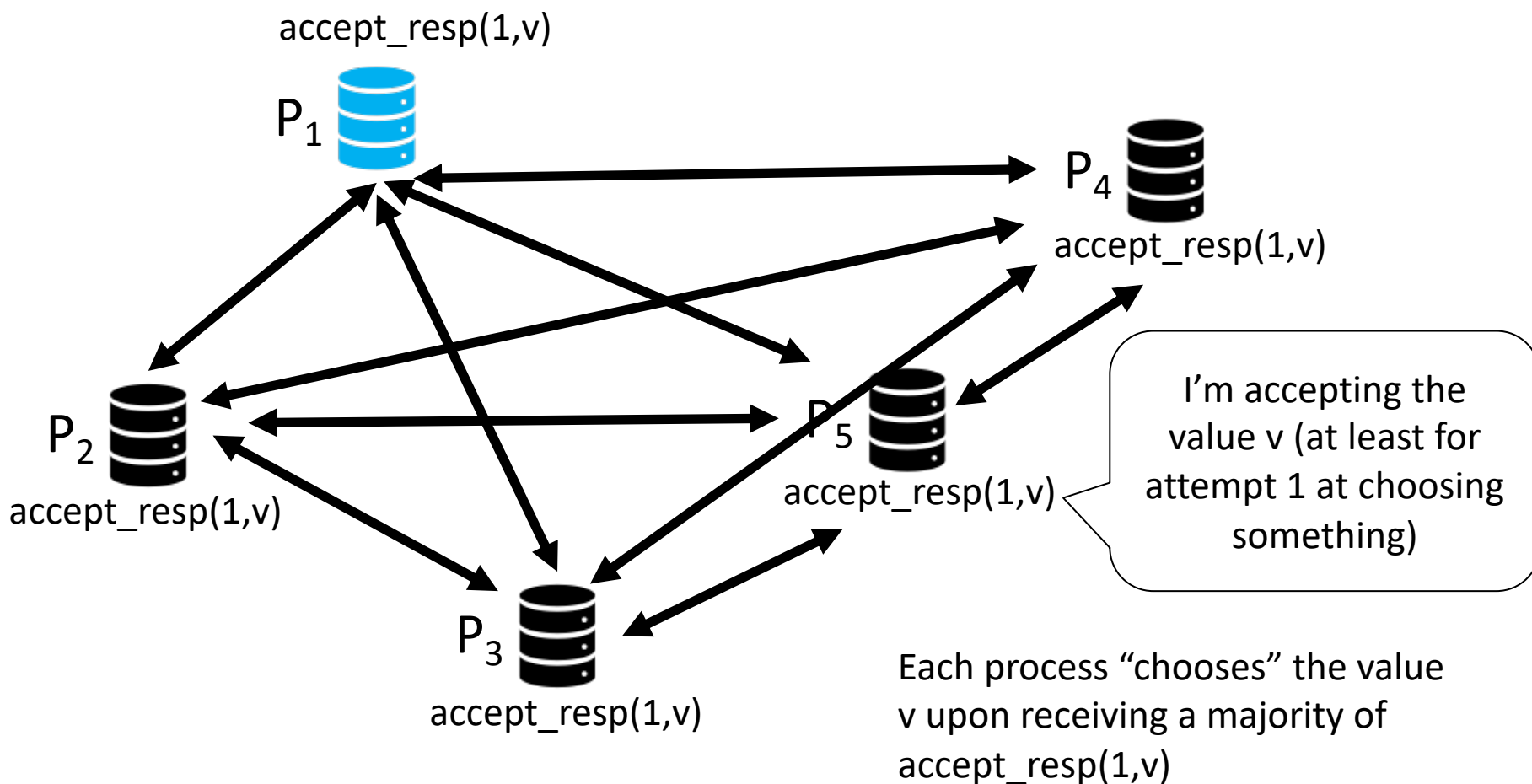
# Normal Case (Option 1): Inform Distinguished Learner

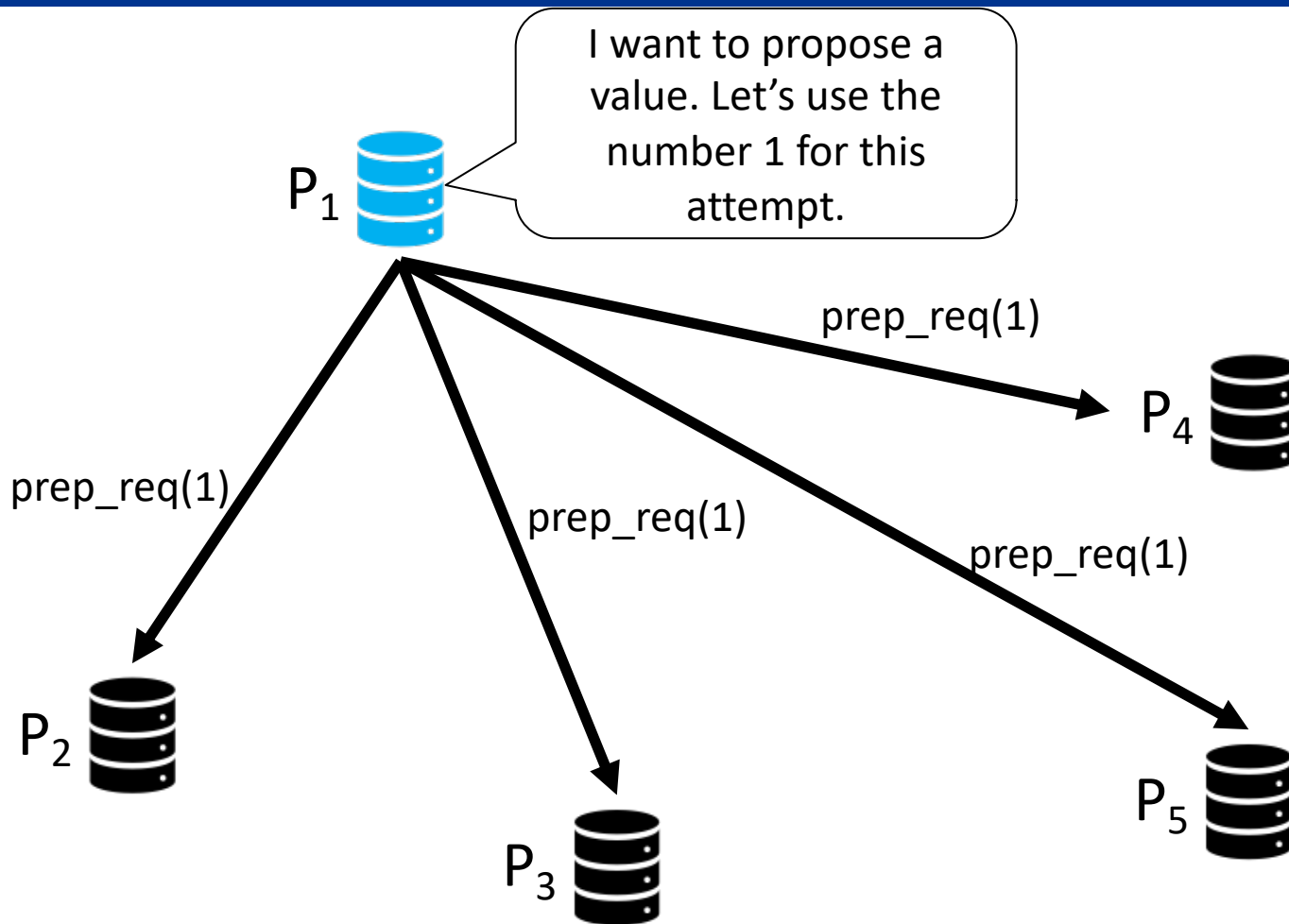# Normal Case (Option 1): Inform Distinguished Learner

# Normal Case (Option 2): Inform all Learners



accept_resp(1,v)

P$_1$

P$_4$

accept_resp(1,v)

P$_2$

P$_5$

accept_resp(1,v)

accept_resp(1,v)

I'm accepting the value v (at least for attempt 1 at choosing something)

P$_3$

accept_resp(1,v)

Each process "chooses" the value v upon receiving a majority of accept_resp(1,v)

# Leader Failure Case: Prepare Response

# Leader Failure Case: Accept Request



P$_1$

Okay, no value could have been accepted, so let's use my value v.

**P$_1$ fails in the middle of sending its accept req!**

P$_4$

accept_req(1,v)

accept_req(1,v)

P$_2$

P$_3$

P$_5$

# Leader Failure Case: Prepare Request (Attempt 2)

$P_2$ and $P_3$ knew about previous proposal with value v, so $P_2$'s proposal is constrained to value v

$P_1$

$P_4$

prep_resp(2,null,0)

$P_2$

prep_resp(2,null,0)

$P_5$

prep_resp(2,v,1)    $P_3$

$P_1$

All processes accept value v in attempt 2 (and can then learn that it was chosen via accept_resp+choose)

$P_4$

accept_req(2,v)

$P_2$

accept_req(2,v)

$P_5$

accept_req(2,v)

$P_3$

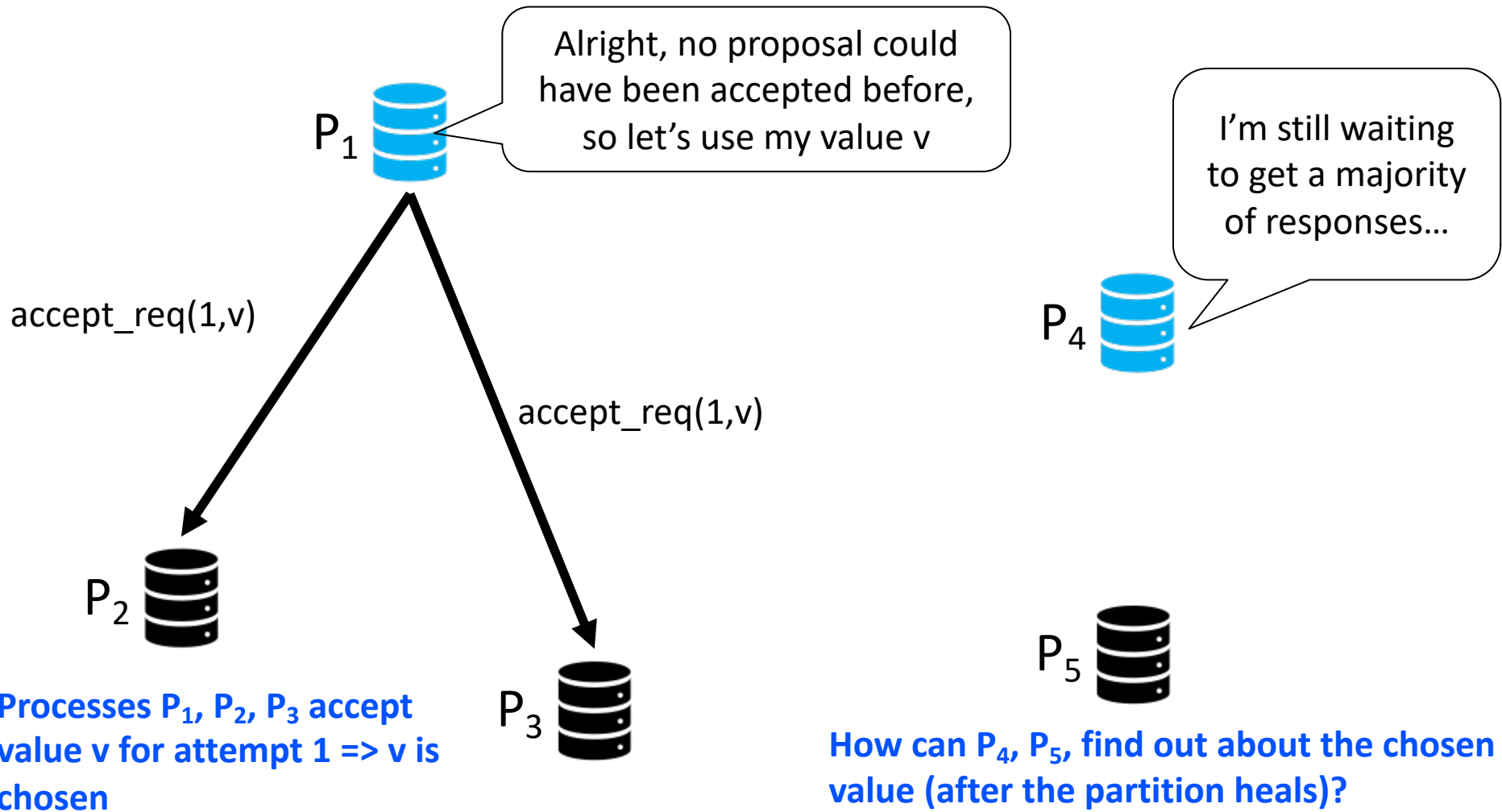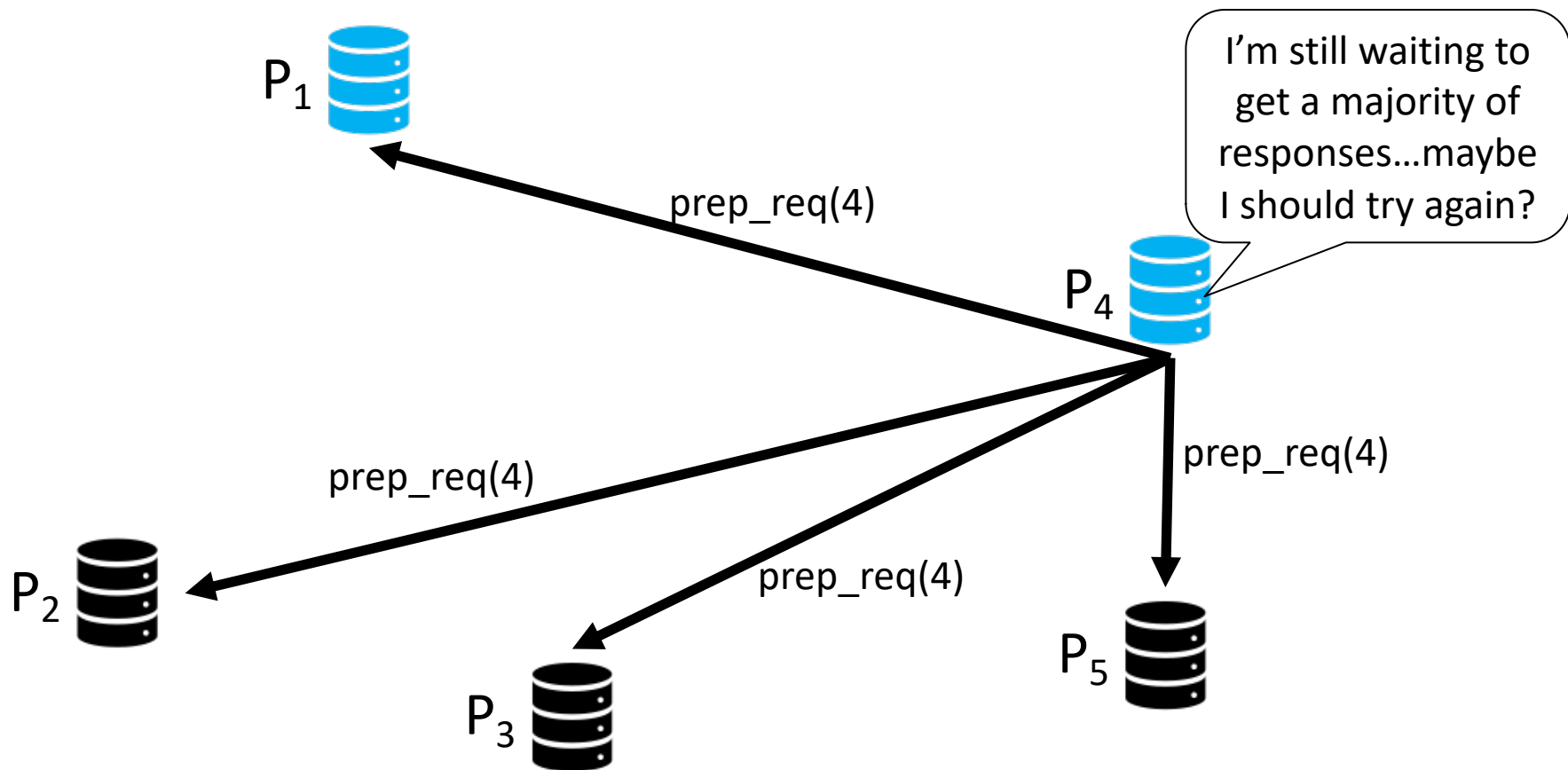# Simple Partition Case: Prepare Request

# Simple Partition Case: Prepare Response

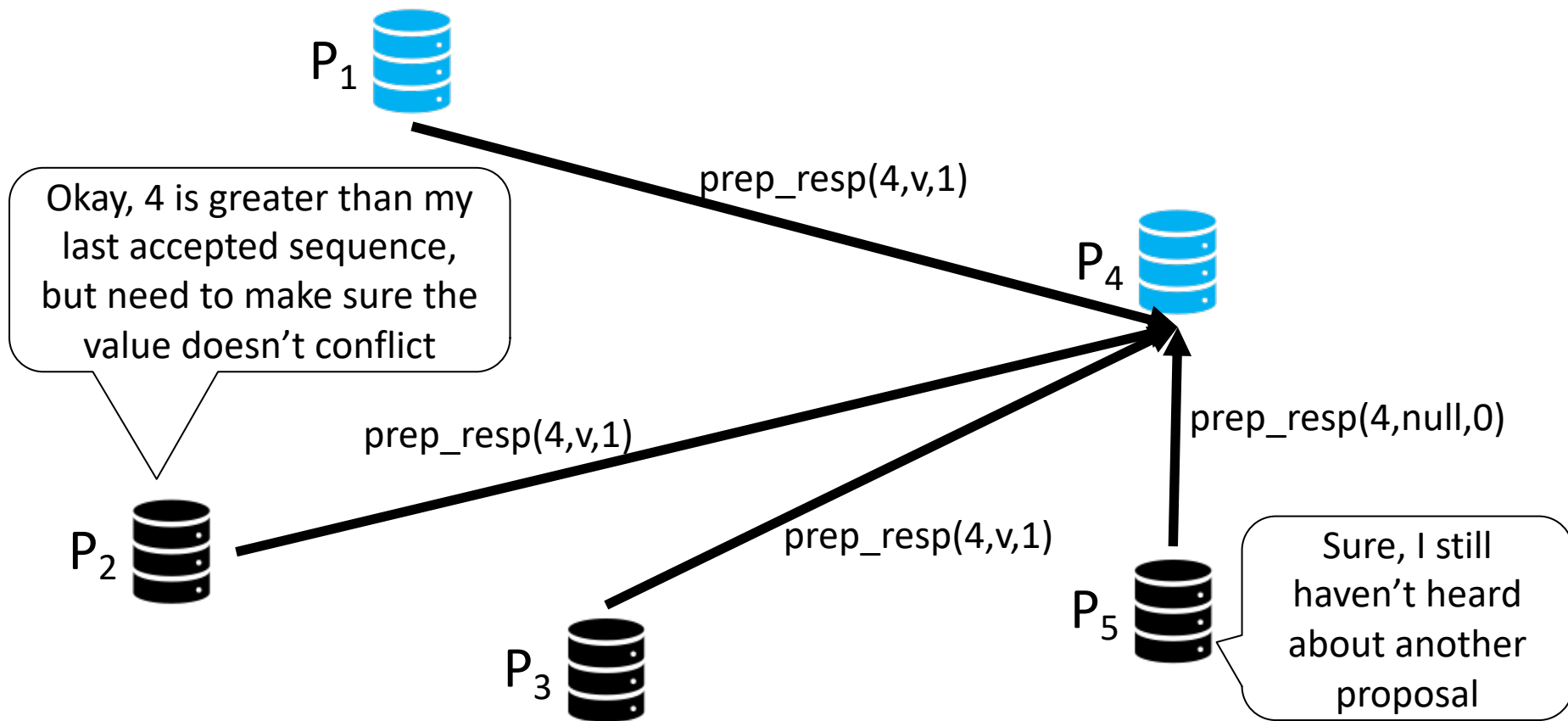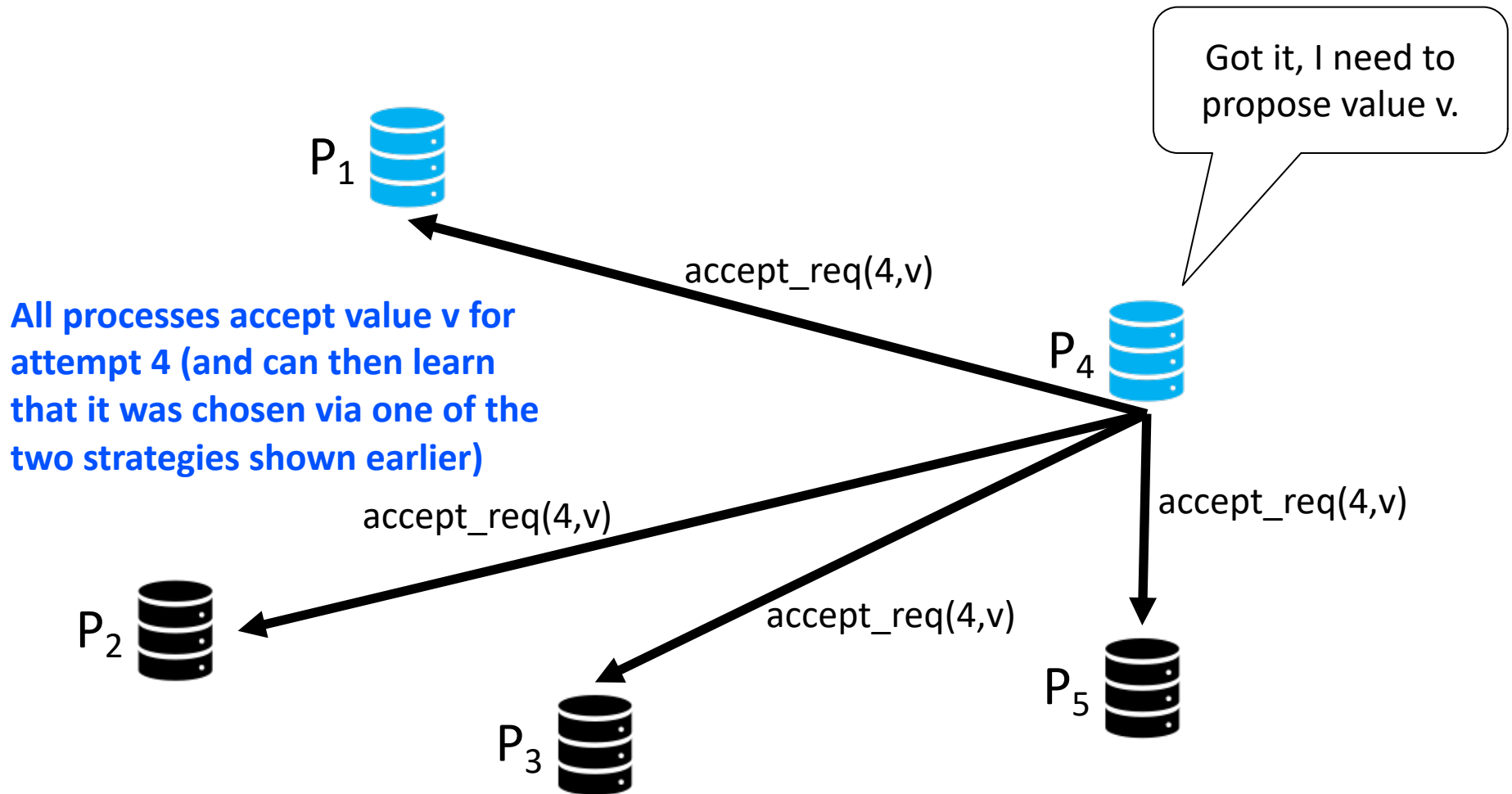# Simple Partition Case: Partition Repair (one option)

# Simple Partition Case: Partition Repair (one option)

# Simple Partition Case: Partition Repair (one option)

# From Consensus to SMR (Multi-Paxos)

- Execute a sequence of separate Paxos instances

- Value chosen in $i^{th}$ consensus instance is the $i^{th}$ command executed

- Optimization: stable leader only needs to execute "prepare" phase once

# Proposer Actions (Prepare phase)

- **Send** prepare_request(n,i)
  - n : sequence/"view" number
  - i : ordinal of last (consecutive) consensus instance for which the proposer knows a *chosen* value
- **Wait** for majority
  prepare_response(n,{$(j,v_j,m_j)$, $(k,v_k,m_k)$, ...})
  - $(j,v_j,m)$ : (ordinal, value, view) for each ordinal j > i for which the acceptor has previously accepted a value $v_j$ in some view $m_j$

# Proposer Actions (Accept phase)

- **Propose all constrained updates:**
  - *For each* ordinal *j* s.t. some $v_j$ was received, **send** accept_request(n,j,$v_j$) where $v_j$ is the value associated with highest received $m_j$
- **Fill all holes (via no-ops):**
  - For each ordinal $j < j_{max}$ s.t. no $v_j$ was received, send accept_request(n,j,no-op)

- **Proceed with unconstrained instances**
  - Incoming client requests can be assigned ordinals starting with $j_{max}+1$

# Acceptor Actions

- **Upon receiving** prepare_request(n,i)
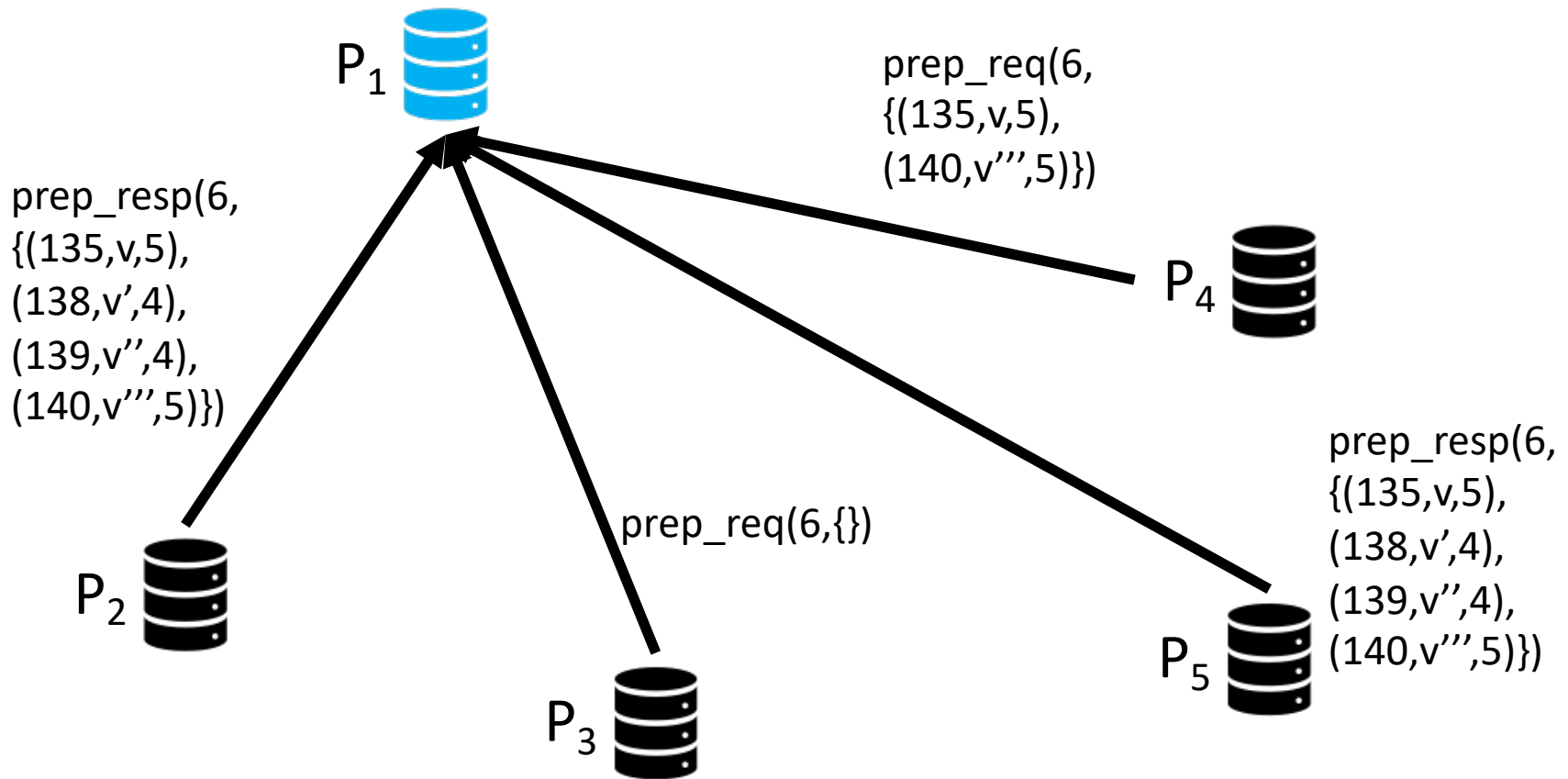  - *If* previously responded to prepare_request(m,*) s.t. m > n, do nothing (or "inform proposer" as performance optimization)
  - *Else* construct constraining update list and send prepare_response(n,{(j,$v_j$,$m_j$), (k,$v_k$,$m_k$), ...})

- **Upon receiving** accept_request(n,j,$v_j$)
  - *If* previously responded to prepare_request(m,*) s.t. m > n, do nothing (or "inform proposer" as performance optimization)
  - *Else* **accept** value $v_j$ for ordinal j (and inform learners)

# Multi-Paxos Example (from paper)



P₁

prep_req(6,
{(135,v,5),
(140,v''',5)})

prep_resp(6,
{(135,v,5),
(138,v',4),
(139,v'',4),
(140,v''',5)})

P₄

prep_resp(6,
{(135,v,5),
(138,v',4),
(139,v'',4),
(140,v''',5)})

P₂

prep_req(6,{})

P₅

P₃

# Multi-Paxos Example (from paper)



accept_req(6, 135,v)

accept_req(6, 135,v)

accept_req(6, 135,v)

accept_req(6, 135,v)

P1

P2

P3

P4

P5

# Multi-Paxos Example (from paper)



P₁

accept_req(6,136,no-op)

accept_req(6, 136, no-op)

P₄

accept_resp(6,135,v)

accept_req(6, 136,no-op)

accept_req(6, 136, no-op)

P₂

P₃

P₅

# Multi-Paxos Example (from paper)



P₁

accept_req(6,137,no-op)

accept_req(6, 137, no-op)

P₄

accept_resp(6,136,no-op)

accept_req(6,137,no-op)

accept_req(6, 137, no-op)

P₂

P₃

P₅

# Multi-Paxos Example (from paper)



P$_1$

accept_req(6,140,v''')

accept_req(6, 140,v''')

P$_4$

accept_resp(6,137,no-op)

accept_req(6,140,v''')

accept_req(6, 140,v''')

P$_2$

P$_3$

P$_5$

# Multi-Paxos Example (from paper)



P₁

accept_req(6,141,X)

accept_req(6, 140,X)

P₄

accept_resp(6,140,v''')

accept_req(6,141,X)

accept_req(6, 140,X)

P₂

P₃

P₅

# Summary

- Introduced the Paxos consensus algorithm (simplified from original specification)

- Key idea behind preventing inconsistency is simple and intuitive

- Implementation and ensuring liveness (especially ensuring that replicas can actually **execute** ordered events) requires non-trivial extensions

  – Often dismissed as "engineering details"