# Distributed Systems } Paxos

But First... Administrivia

- Midterm: Next Week, In Class
    - Covers Everything Up To Today
    - Open Book
        ↳ Papers, Notes, Etc.
        → Electronic Devices
            ↳ Can Use Ipad Etc.
            But
                - Must Disconnect From Internet: Download Ahead Of Time
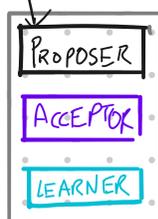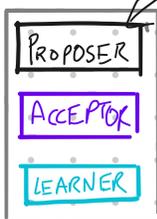                - No Chat Bots, Etc.

— I uploaded last year's midterm to the
website if you wanted an idea of
what to expect.

A brief digression ———>

# Paxos

- Why? Reinforce ideas around correctness

- General form: 3 types of nodes.
  SYNOD: Log with 1 slot/index

| PROPOSER | ACCEPTOR | LEARNER |
| PROPOSER | ACCEPTOR | LEARNER |
| PROPOSER | ACCEPTOR | LEARNER |

① PROPOSE VALUES:

Suggest command
to set index to

$f+1$

② ACCEPT VALUES:
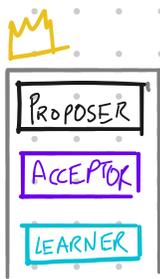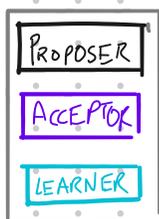
CHOOSE WHAT
VALUE TO ACCEPT
A COMMIT

$2f+1$
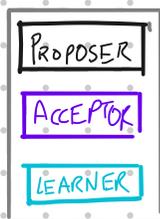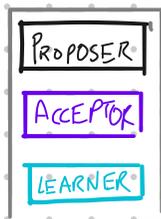
③ LEARN VALUES:

EXECUTE COMMAND

$f+1$

WHY SEPARATE THIS OUT?

RECIPE FOR CONSTRUCTING RSM PROTOCOLS

↳ DIFFERENT TYPES OF PROPOSERS/ACCEPTORS/ LEARNERS

RETURN TO THIS AT THE END OF CLASS.

MULTIPAXOS /CHUBBY : OUR SETTING

| PROPOSER |
|----------|
| ACCEPTOR |
| LEARNER |

| PROPOSER |
|----------|
| ACCEPTOR |
| LEARNER |

| PROPOSER |
|----------|
| ACCEPTOR |
| LEARNER |

Inactive*

| PROPOSER |
|----------|
| ACCEPTOR |
| LEARNER |

| PROPOSER |
|----------|
| ACCEPTOR |
| LEARNER |

| PROPOSER |
|----------|
| ACCEPTOR |
| LEARNER |

# SYNOD: THE ALGORITHM FOR A SINGLE SLOT
## (WE LOOKED AT THIS LAST CLASS)

GOES IN ROUNDS, EACH ROUND DENOTED BY A "PROPOSAL #"

REQUIREMENTS

COMMIT

- CHOOSE SOME VALUE, EVENTUALLY (Avoids trivial solutions)

P1  ↳ ACCEPTOR ACCEPTS FIRST VALUE
                ↳ Not the same as choosing!

P2 - ONCE A VALUE V IS CHOSEN, NO OTHER VALUE

WILL BE CHOSEN

↳ P2a. IF VALUE $V$ IS CHOSEN AS PROPOSAL $N$, THEN ANY PROPOSAL WITH # $> N$ ACCEPTED BY <u>AN ACCEPTOR</u> MUST PROPOSE $V$.

[NOTE, INCLUDES VALUES ACCEPTED BY ACCEPTORS WHO HAVE MISSED ALL PROPOSALS/ MESSAGES So FAR]

P2b. IF VALUE $V$ IS CHOSEN AS PROPOSAL $N$, ~~THEN~~

~~ANY PROPOSAL WITH # $> N$ ACCEPTED BY AN ACCEPTOR~~ ~~MUST HAVE VALUE $V$~~ ANY HIGHER # PROPOSAL ISSUED BY A PROPOSER HAS VALUE $V$.

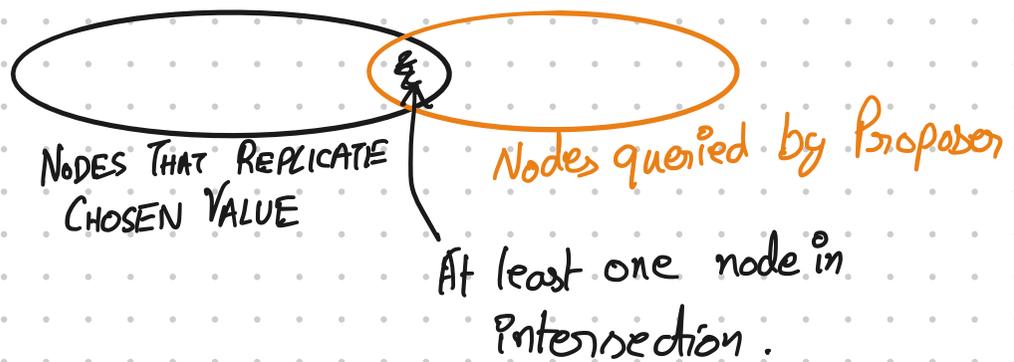[Note, must be true even for proposals issued by a different proposer in the future. ]

Observe - P2b $\Rightarrow$ P2a

- For P2b to hold proposers need to be able to find any chosen (committed value $V$
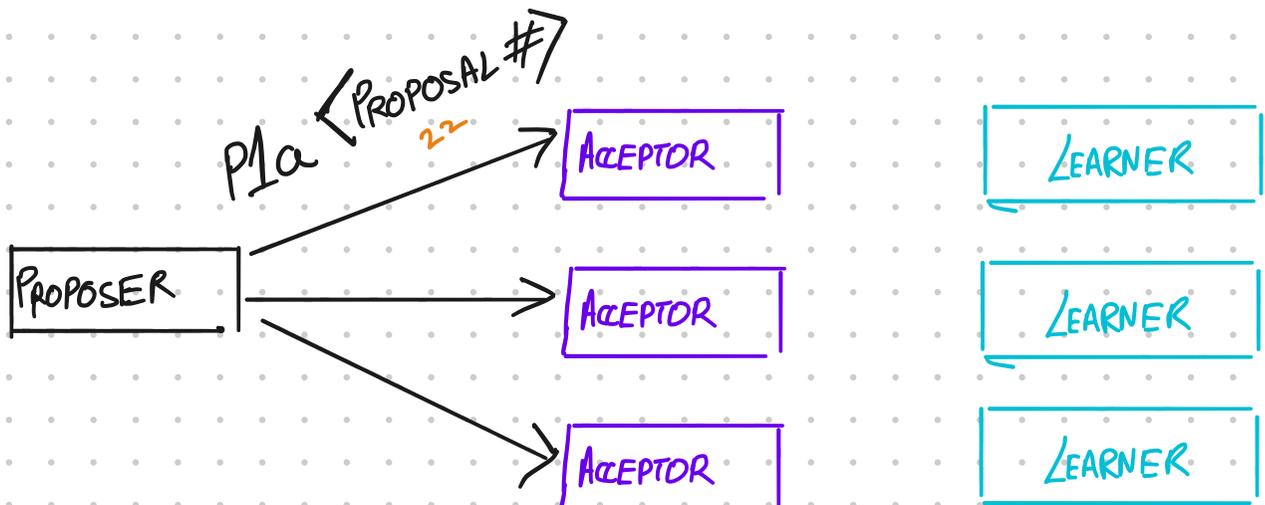
- CAN USE QUORUM INTERSECTION FOR THIS

(P2C)

REQ: ① CHOSEN VALUES MUST BE ACCEPTED
BY A QUORUM OF ACCEPTORS

(COMMITTED → over CHOSEN)

② PROPOSER MUST CHECK WITH A
QUORUM BEFORE PROPOSING A VALUE

NODES THAT REPLICATE    Nodes queried by Proposer
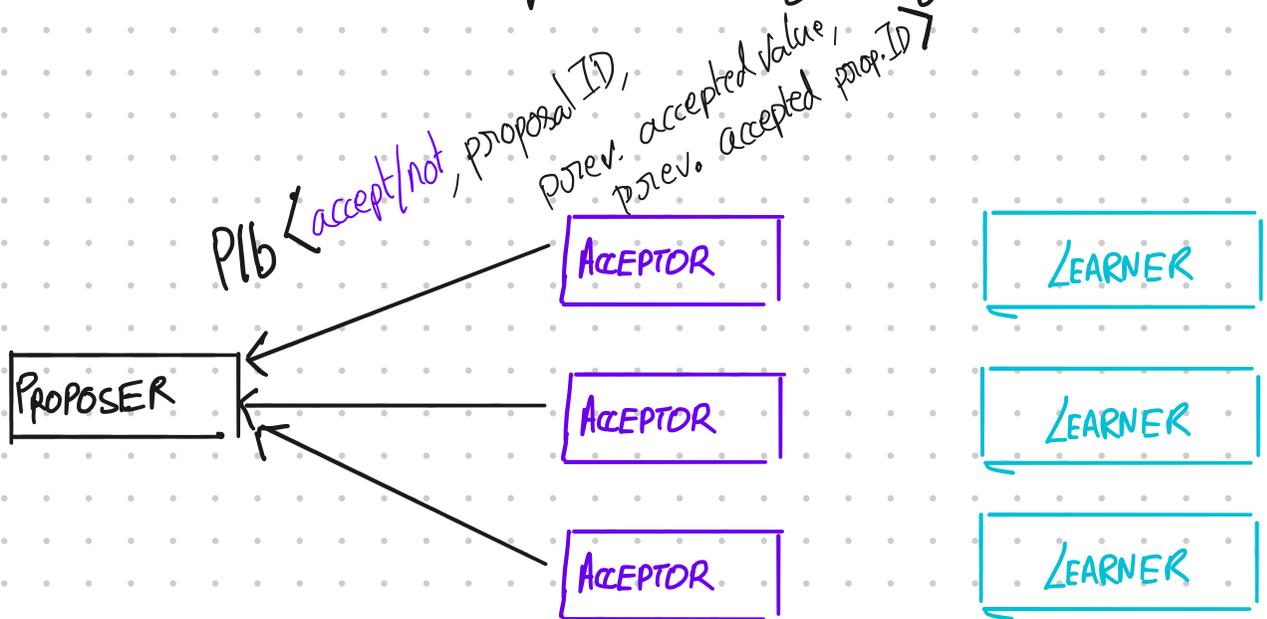CHOSEN VALUE

At least one node in
intersection.

$P2a/P2b \Rightarrow$ If a value $V$ was committed w/
proposal # $P$, all accepted proposals w/
prop# $p+1$ or larger contain $V$

$\hookrightarrow$ Value associated with largest proposal
# is the one likely committed

P1a { PROPOSAL #
       22

PROPOSER → ACCEPTOR          LEARNER

→ ACCEPTOR          LEARNER

→ ACCEPTOR          LEARNER

Goals ⓐ Query to find any previously accepted values.

ⓑ Prevent any pending proposals (with lower ID) from making progress

$P/b \langle$ *accept/not*, proposal ID, prev. accepted value, prev. accepted prop.ID $\rangle$

```
              ┌──────────┐          ┌──────────┐
              │ ACCEPTOR │          │ LEARNER  │
              └──────────┘          └──────────┘
┌──────────┐  ┌──────────┐          ┌──────────┐
│ PROPOSER │◄─│ ACCEPTOR │          │ LEARNER  │
└──────────┘  └──────────┘          └──────────┘
              ┌──────────┐          ┌──────────┐
              │ ACCEPTOR │          │ LEARNER  │
              └──────────┘          └──────────┘
```

∴ FALSE, 24, ... ─     ← 22 is too low

─ True, 22, ⊥, ─     ← Nothing accepted

─ True, 22, X, 2 ⇐ X accepted as prop. ID 2

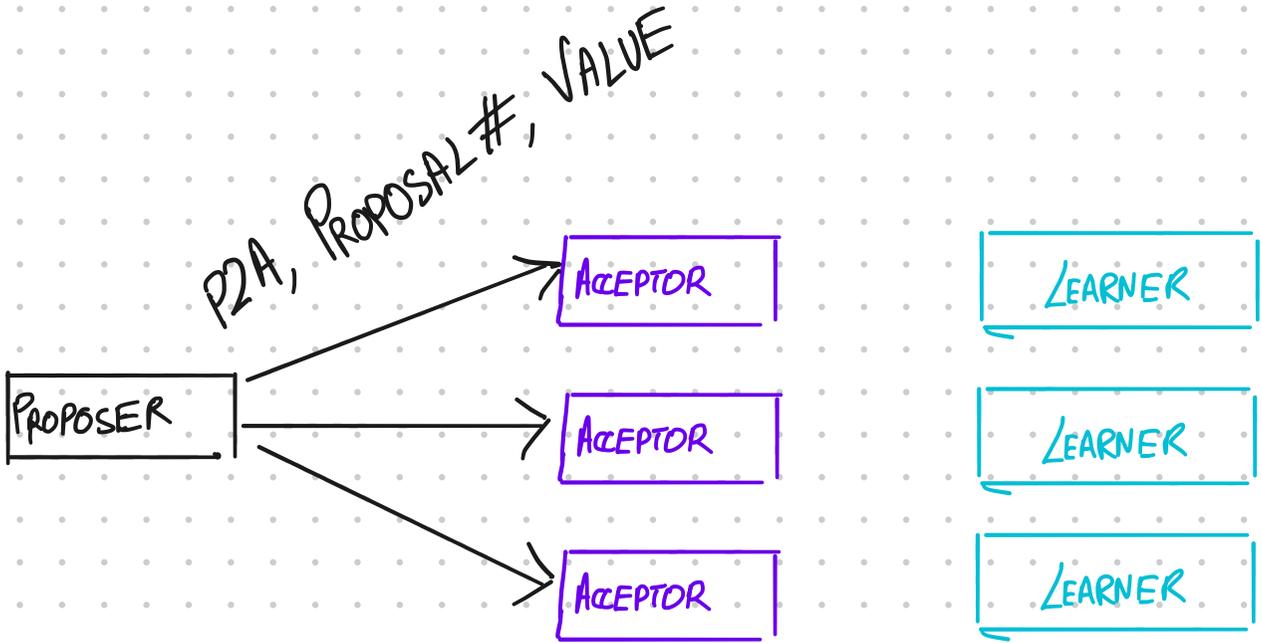n=3                    Value to propose        Possibly Committed
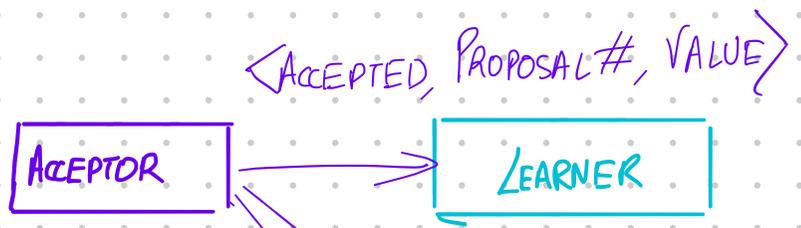
( ⊥, ⊥ )                     Any                    ⊗

( X, 7 )        ⊥             X

$(x, 2)$    $(y, 3)$

$y$

$\{y, \perp\}$

P2A, PROPOSAL#, VALUE

PROPOSER

ACCEPTOR

ACCEPTOR

ACCEPTOR

LEARNER

LEARNER

LEARNER

When should the ACCEPTOR not accept?

⟨ACCEPTED, PROPOSAL#, VALUE⟩

ACCEPTOR

LEARNER

PROPOSER

ACCEPTOR

ACCEPTOR

LEARNER

LEARNER

- When should *learners* execute/act on value?

# Why leader?

P1a 22

PROPOSER

PROMISE

ACCEPTOR

ACCEPTOR

ACCEPTOR

LEARNER

LEARNER

LEARNER

PROPOSE

PROPOSER

PROPOSE

ACCEPTOR

ACCEPTOR

ACCEPTOR

LEARNER

LEARNER

LEARNER

P1a 23

# Duelling proposers are a problem!

## Putting it all together

So far : Focus on a single log idx. But what about
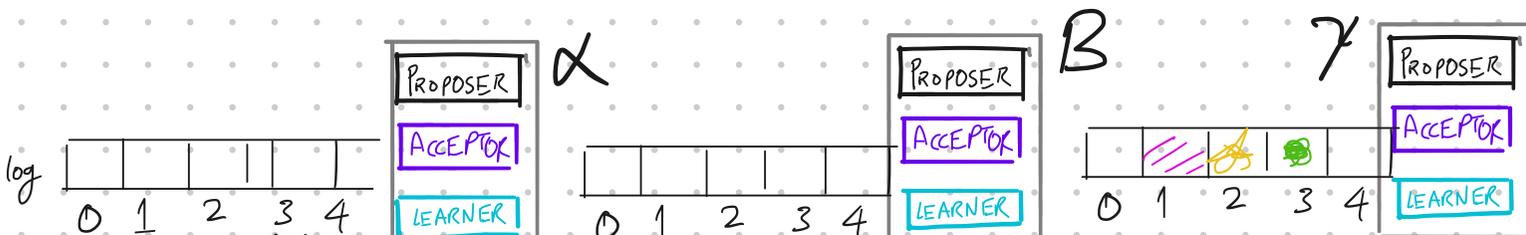. RSMs?     "Fast Paxos" Lamport '05

⟨command, proposal #⟩

log

| 0 | 1 | 2 | 3 | 4 |
View ID/Ballot

**PROPOSER** α
ACCEPTOR
LEARNER

| 0 | 1 | 2 | 3 | 4 |

**PROPOSER** B
ACCEPTOR
LEARNER

**PROPOSER** γ

| 0 | 1 | 2 | 3 | 4 |

ACCEPTOR
LEARNER

Log index

P1A, view ID* [0...∞]

T
I
M
E  ↓

Note : Only require
that proposal #
are s.t. it
v1 < v2 then
prop# in v1 <
        prop# in v2

log

| 0 | 1 | 2 | 3 | 4 |

**PROPOSER** α
ACCEPTOR
LEARNER

| 0 | 1 | 2 | 3 | 4 |

**PROPOSER** B
ACCEPTOR
LEARNER

B

| 0 | 1 | 2 | 3 | 4 |

γ

**PROPOSER**
ACCEPTOR
LEARNER

View ID/Ballot

T
I
M
E
↓

P1A, view ID,* [0∞∞]

d

True if this
is highest view ID!

⟨P1b, [☐☐] , view⟩

| | 0 | 1 | 2 | 3 |
|---|---|---|---|---|

Wait for P1b from a quorum.
Merge logs

| | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| α | b,1 | (x,1) | y,1 | z,1 |
| γ | (⊥) | a,2 | y,1 | z,1 |
| ML : | b,1 | a,2 | y,1 | z,1 |

log

| | 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|

0  1  2  3  4
View ID/Ballot

| PROPOSER |
|---|
| ACCEPTOR |
| LEARNER |

α

| | 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|

| PROPOSER |
|---|
| ACCEPTOR |
| LEARNER |

B

d

γ

| | 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|

| PROPOSER |
|---|
| ACCEPTOR |
| LEARNER |

P1A,  4  , [0∞∞]

TIME ↓

$\langle P1b, \cdots \rangle$

— Merge

P2a, 4;

    [b, a, y, z]

Update Log

| b,4 | a,4 | y,4 | z,4 | |
|---|---|---|---|---|

P2b

P2b

**Compare to Raft!!!**

Count for each log index.

Execute when P2b from quorum

P2A, 4, 4→c

c

- - -

Electing a leader?

    - Only requirement is Election Safety;
      that is at most one leader

    - Any protocol that meets this requirement

Suffices

Changing Configuration

- Hard problem, for the same reason as last week : need to avoid quorums of nodes that are not up-to-date

Vertical Paxos.

What does this generality help with

- Disk Paxos
  ↳ Acceptors and learners are networked disks
      ↳ Read, Write

- Mencius
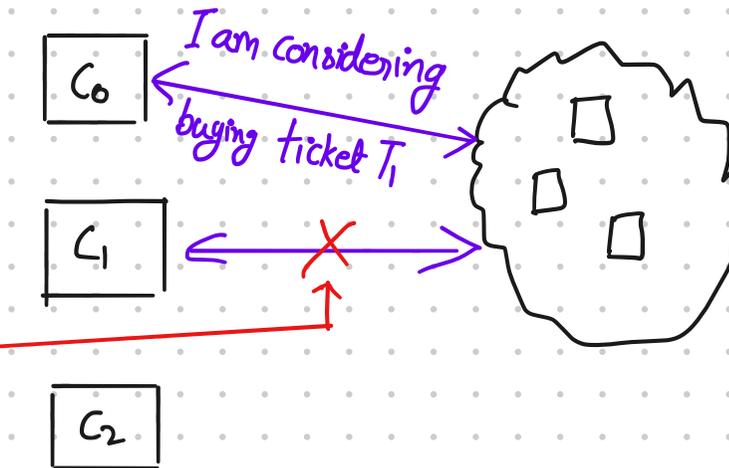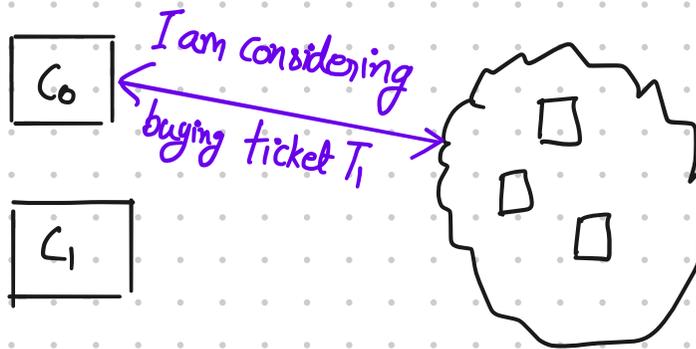  ↳ Multiple leaders for multiple datacenters
  ⇒ Later improved by EPaxos
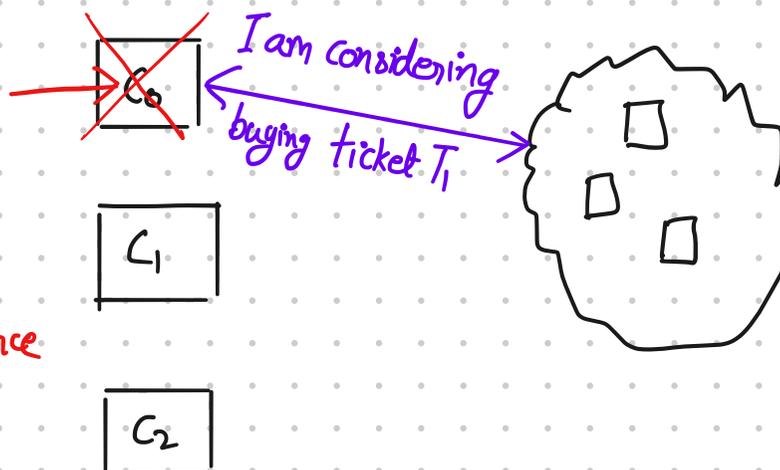
. . .

Taking a step back

- RSMs

- Quorum intersection

↳ Leader safety

↳ State machine safety

↳ Leader safety

↳ State machine safety

# Leader Leases / Leases

## Lease

$C_0$ ← I am considering buying ticket $T_1$ →

$C_1$

$C_2$

---

$C_0$ ← I am considering buying ticket $T_1$ →

$C_1$ ⇐ ✗ →

**Don't offer $C_1$ the same ticket**

$C_2$

---

**Make sure failure at $C_0$ cannot cause $T_1$ to be unsold forever [Resource Leak]**

~~$C_0$~~ ← I am considering buying ticket $T_1$ →

$C_1$

$C_2$

A common problem
- Locks (mutex)
- Resources (memory, etc.)
- ...

Leases [Gray & Cheriton '89]

$C_0$ — I am considering buying ticket $T_1$ → (cloud)

You have 1 minute to decide & buy

$C_1$

LEASE

$C_2$

---

$C_0$

need more time

Have another minute

Renewal

$C_1$

$C_2$

---

$C_0$

$C_1$

Not shown $T_1$

until lease
expires.

$\boxed{C_2}$

Observe: Raft HB/Leader Election is a Lease Mechanism



T=2

A (crown) —— AER/HB ——> B  T=2
         ——> (cyan) <——
A —— AER/HB ——> C  T=2

(cyan, from B to A) Grant: Promise to not
start a new election
for ____



T=2  A (crown)

B

         ↑ RV

Rv ← C
      T=3        ← Lease expired

---

$\boxed{\text{Leader Leases}}$

Goal:

- Stable leadership
- Availability.