G

# TEE + Raft →
# Confidential Replicated State Machines

Dzmitry Huba

# Introduction

# Confidential services landscape

## Scenarios

**Lock service**
Policy bookkeeping and enforcement

**Key value store**
Transactional data store for private data

**Key provisioning service**
Key management and distribution for trusted workloads

**Private data processing**
Analytics or learning over private data

## Properties

**Stateful**
Uses private state throughout the lifetime

**Long-running**
Lifetime measures in days

**Strongly consistent**
Necessary to guarantee privacy

**Multiple releases**
Partial results are released multiple times throughout the lifetime

## Challenges

**Scalability**
Ranges from small to massive scale

**Failure resilience**
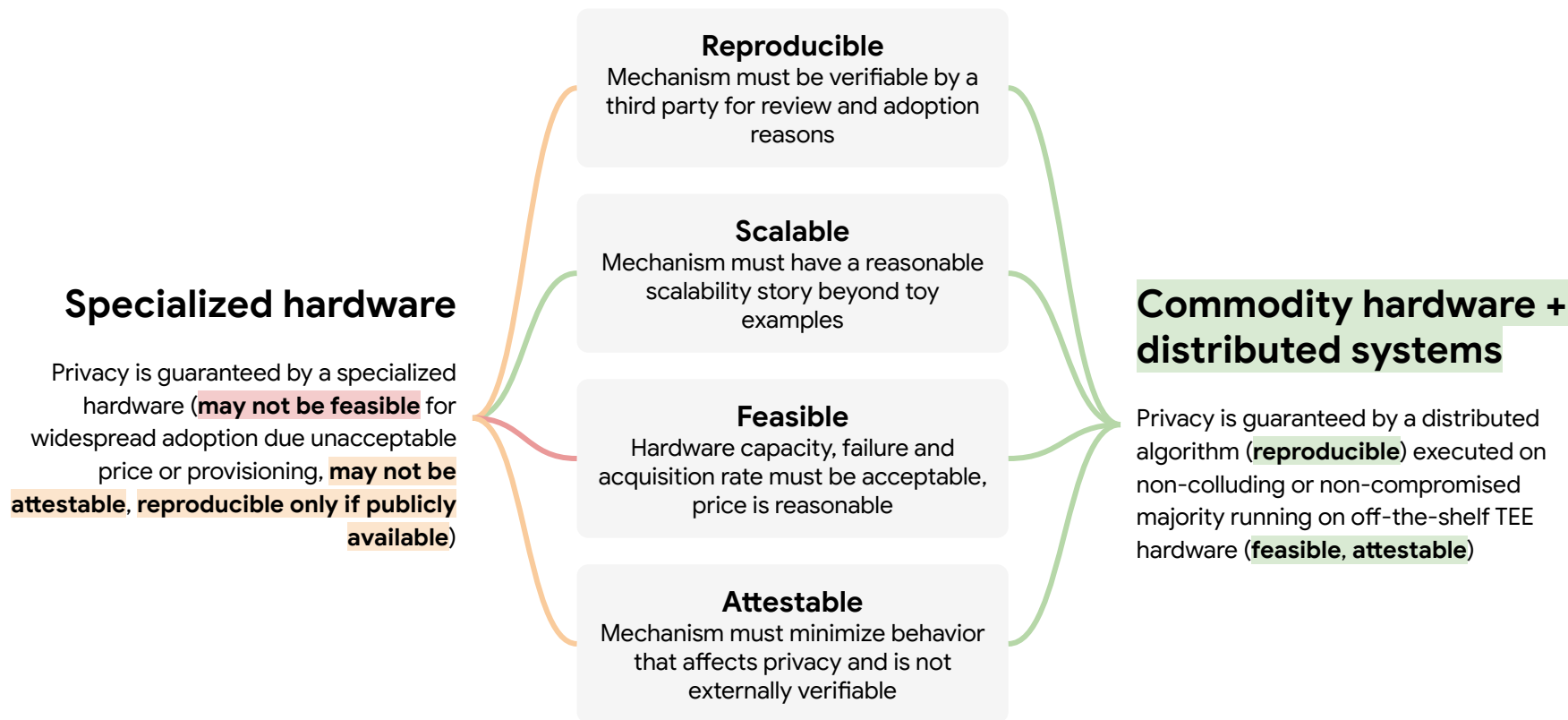Failures on the days long timeframe are likely to happen

**Protection against attacks**
Attack surface grows with complexity

**Understandability**
Critical for review process and overall narrative

# Privacy mechanisms options and desirable properties

## Specialized hardware

Privacy is guaranteed by a specialized hardware (**may not be feasible** for widespread adoption due unacceptable price or provisioning, **may not be attestable**, **reproducible only if publicly available**)

### Reproducible
Mechanism must be verifiable by a third party for review and adoption reasons

### Scalable
Mechanism must have a reasonable scalability story beyond toy examples

### Feasible
Hardware capacity, failure and acquisition rate must be acceptable, price is reasonable

### Attestable
Mechanism must minimize behavior that affects privacy and is not externally verifiable

## Commodity hardware + distributed systems

Privacy is guaranteed by a distributed algorithm (**reproducible**) executed on non-colluding or non-compromised majority running on off-the-shelf TEE hardware (**feasible**, **attestable**)

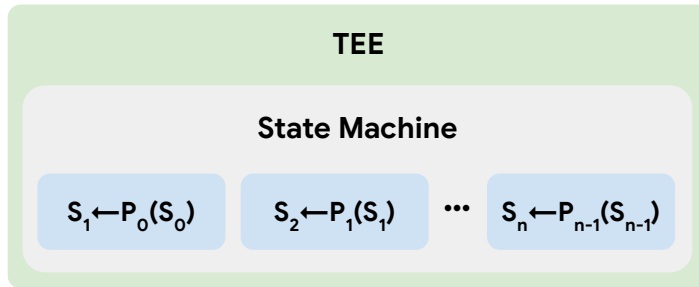Replicated state machines is a general method for implementing a fault-tolerant service

Critical step is choosing an order for the inputs to be processed, can be done with consensus protocol (Raft!)

TEE + Raft enables replicated state machines with rollback protection and externally verifiable behavior
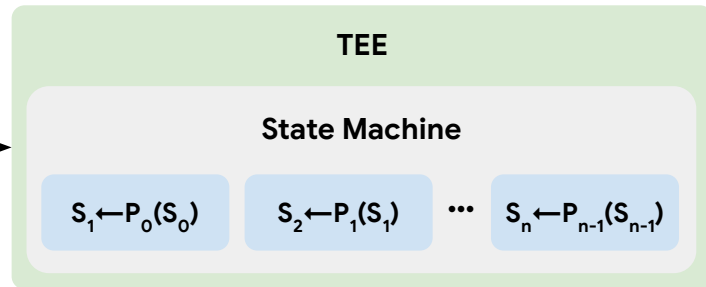
# Physical TEE

✓ **Confidentiality**
**Isolation**
**Verifiability**

**TEE**

**State Machine**

$S_1 \leftarrow P_0(S_0)$    $S_2 \leftarrow P_1(S_1)$    $\cdots$    $S_n \leftarrow P_{n-1}(S_{n-1})$

✗ **Fault tolerance**
**Scalability**

# Replicated TEE

**Confidentiality Isolation Verifiability** ✓

**Fault tolerance Scalability** ✓

**Subject to rollback and forking attacks** ✗

**Replicated state machine on top of Raft consensus protocol**

**TEE**

State Machine

$S_1 \leftarrow P_0(S_0)$   $S_2 \leftarrow P_1(S_1)$   $\cdots$   $S_n \leftarrow P_{n-1}(S_{n-1})$

**TEE**

State Machine

$S_1 \leftarrow P_0(S_0)$   $S_2 \leftarrow P_1(S_1)$   $\cdots$   $S_n \leftarrow P_{n-1}(S_{n-1})$

**TEE**

State Machine

$S_1 \leftarrow P_0(S_0)$   $S_2 \leftarrow P_1(S_1)$   $\cdots$   $S_n \leftarrow P_{n-1}(S_{n-1})$

**Service TEE**

Identity protection

End to end encryption

✓ Confidentiality
Isolation
Verifiability

**TEE**

State Machine

$S_1 \leftarrow P_0(S_0)$   $S_2 \leftarrow P_1(S_1)$  ···  $S_n \leftarrow P_{n-1}(S_{n-1})$

Log integrity protection

Replicated state machine on top of Hardened Raft consensus protocol

**TEE**

State Machine

$S_1 \leftarrow P_0(S_0)$   $S_2 \leftarrow P_1(S_1)$  ···  $S_n \leftarrow P_{n-1}(S_{n-1})$

Mutual attestation

✓ Fault tolerance
Scalability

✓ Protected from rollback and forking attacks

**TEE**

State Machine

$S_1 \leftarrow P_0(S_0)$   $S_2 \leftarrow P_1(S_1)$  ···  $S_n \leftarrow P_{n-1}(S_{n-1})$

Lift TEE guarantees
from process to service level

# Building confidential services

**Key Provisioning Service**

Key → Policy

**Lock Service**

Lock → Metadata, Owner

✓ Powerful primitive, lowers entry barrier for implementation and reasoning

## Service TEE

### Replicated State Machine

$S_1 \leftarrow P_0(S_0)$  $S_2 \leftarrow P_1(S_1)$  ···  $S_n \leftarrow P_{n-1}(S_{n-1})$

✓ Run managed or self-hosted deployment

**Key Value Store**

Key → Value

**Private Data Processing**

Data → Privacy Workflow → State

# Programming model

# Confidential computations are represented through a simple **actor model**

**Commands**

$C_1$  $C_2$  ...  $C_n$

**2** Replicate

**Events**

$E_1$  $E_2$  ...  $E_m$

**3** Apply

**a** Load snapshot

**State**

$S$

**b** Save snapshot

**4** Reply  **1** Propose

**Consumer**

**1** Consumer proposes command for execution

**2** Actor generates event, system replicates event into event log

**3** System commits event after replication to majority, actor applies event to state

**4** Actor replies to consumer

**a** Actor loads state from snapshot when system asks

**b** Actor saves state to snapshot when system asks

# Consumers propose compare and swap operation, once replicated and committed the operation is applied

**a** Load snapshot

| Commands | | Events | | State |
|---|---|---|---|---|
| $C_1$: cas (name, old, new) | **2** Replicate → | $E_1$: cas (name, old, new) | **3** Apply → | S: name → value |

**b** Save snapshot

**4** Reply    **1** Propose

**Consumer**

**1** Req → $C_1$      **3** S[name].cas(old, new)      **a** S ← deserialize(B)

**2** $C_1$ → $E_1$      **4** Resp ← (old, new)      **b** B ← serialize(S)

Actor implementation

# Demo

# Consumers update assigned table cell with optimistic concurrency

|  | Col 1 | Col 2 | Col 3 | Col 4 |
|---|---|---|---|---|
| Row 1 | | | | |
| Row 2 | | **+1** | | |
| ... | ... | ... | ... | ... |
| Row N | | | | |

Consumer 1

Consumer 2

Read row, increment cell value, write if row hasn't changed

...

Consumer M

# Consumers update 16 bits of a named 64 bits counter with a compare and swap operation

|  | 63 - 48 | 47 - 32 | 31 - 16 | 15 - 0 |
|---|---|---|---|---|
| "Row 1" |  |  |  |  |
| "Row 2" |  | +1 |  |  |
| ... | ... | ... | ... | ... |
| "Row N" |  |  |  |  |

Consumer 1

Consumer 2

Consumer M

Read counter, increment bit subset, write if counter hasn't changed

Raft

Best source is
the paper and the visualization

# Raft in a nutshell

## Focus

Makes design decisions based on **understandability**. Uses **problem decomposition** and **state space minimization** (handle multiple problems with single mechanism, eliminate special cases, minimize non-determinism, etc.).
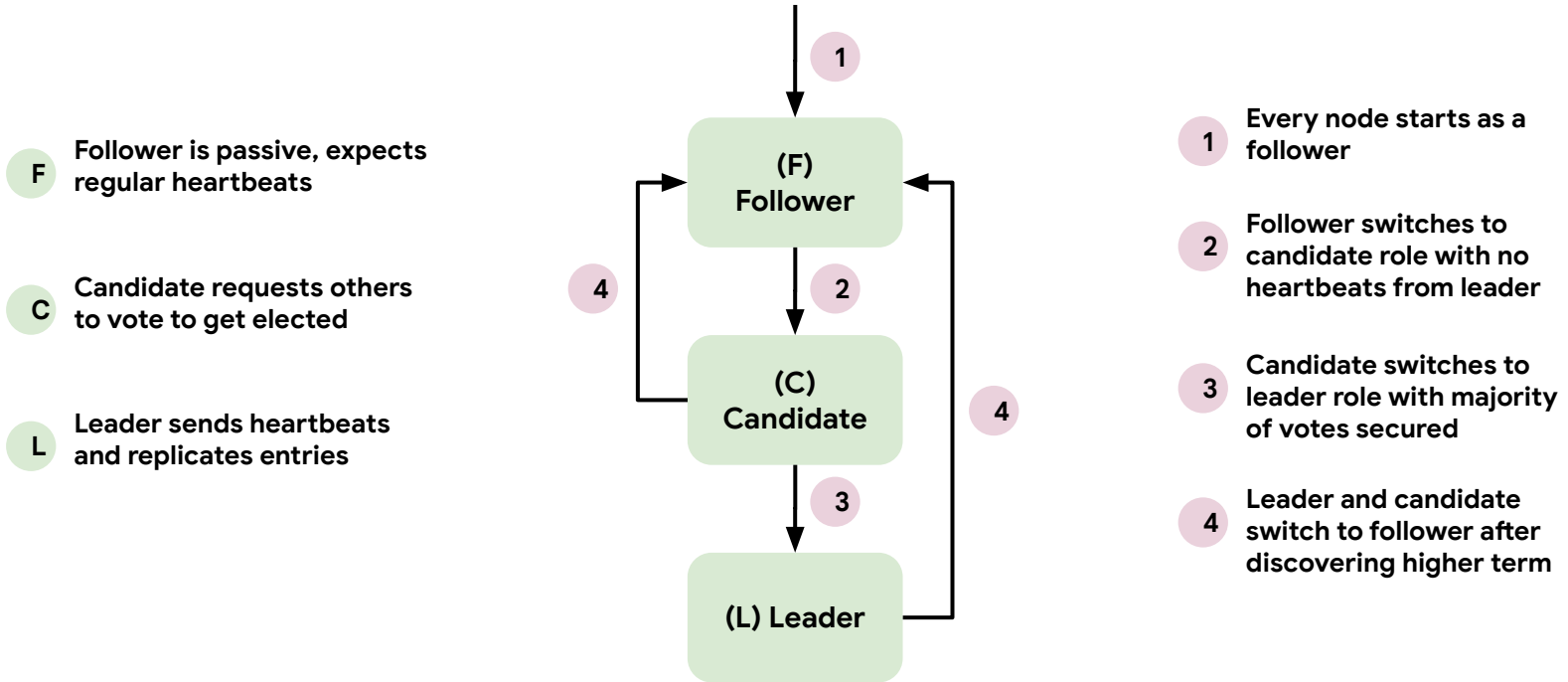
## Decomposition

Consists of **leader election** (select one server to act as leader, detect crashes, choose new leader), **log replication** (leader accepts commands from clients, appends to its log, replicates its log to other servers and overwrites inconsistencies), **safety** (keep logs consistent, only servers with up-to-date logs can become leader).
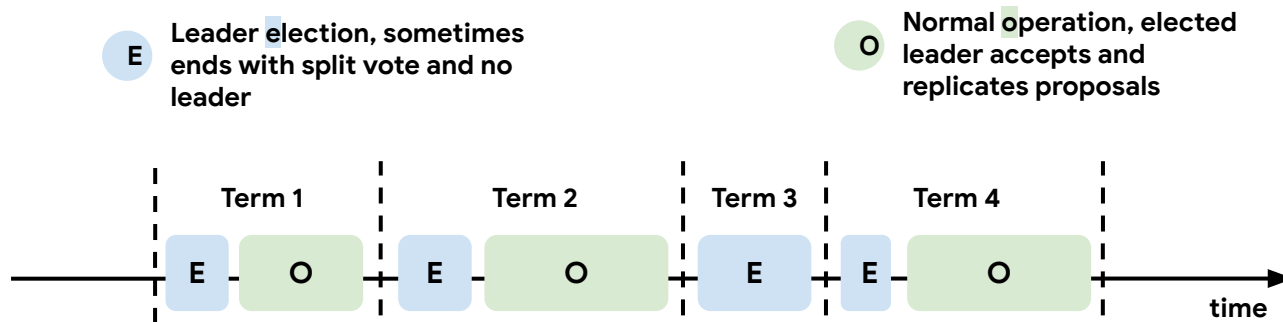
# Replicated state machine



Steps:

1. Consumer issues proposals to the leader Raft node

2. Consensus module appends proposal entry to the log and replicates log to the follower Raft nodes

3. The committed log entry (replicated to the majority) is applied to the state machine
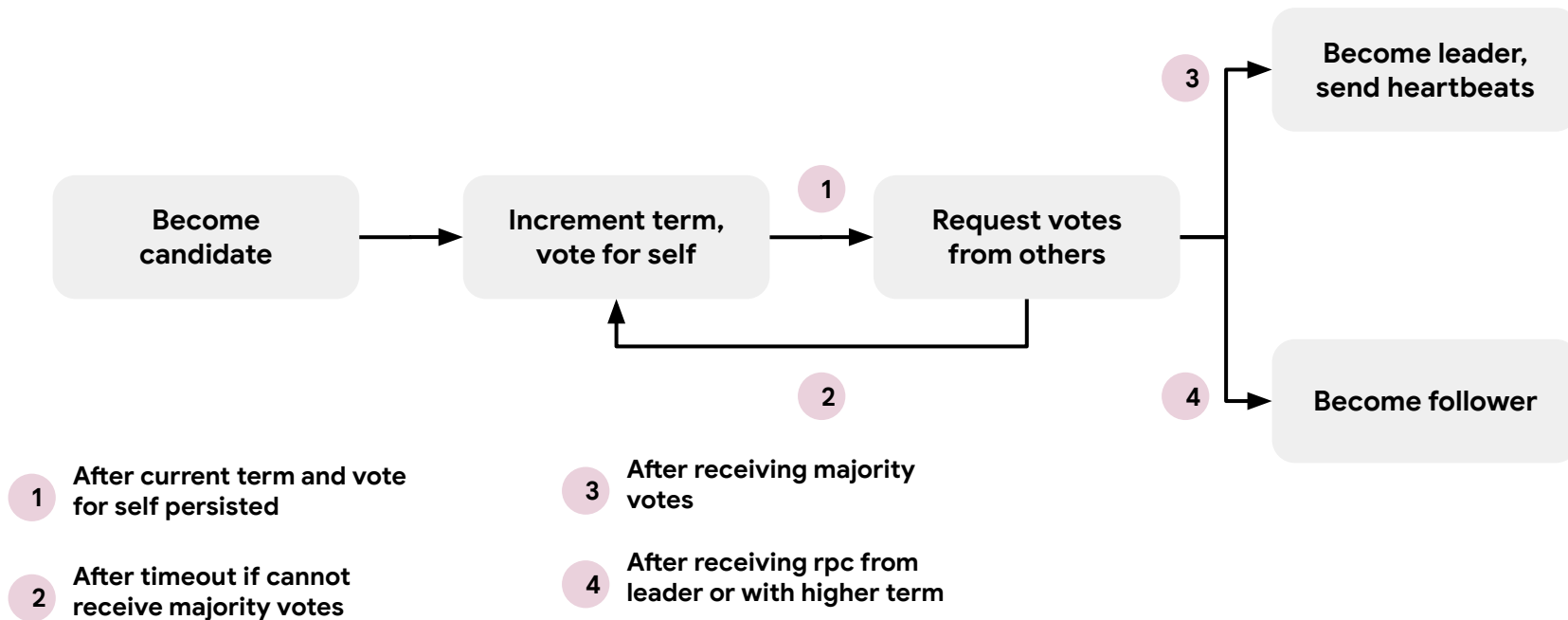
4. The result is sent to the client

# Node states



**F** Follower is passive, expects regular heartbeats

**C** Candidate requests others to vote to get elected

**L** Leader sends heartbeats and replicates entries

**(F) Follower**

**(C) Candidate**

**(L) Leader**

1

2

3

4

4

**1** Every node starts as a follower

**2** Follower switches to candidate role with no heartbeats from leader

**3** Candidate switches to leader role with majority of votes secured

**4** Leader and candidate switch to follower after discovering higher term

Slide is borrowed from author's presentation, restyled

# Terms



**E** Leader election, sometimes ends with split vote and no leader

**O** Normal operation, elected leader accepts and replicates proposals

Term 1 · Term 2 · Term 3 · Term 4

E — O | E — O | E | E — O → time

**1** At most one leader per term, some terms have no leader due to failed election

**2** Each node maintains current term value, exchanged with peers, used to identify obsolete information
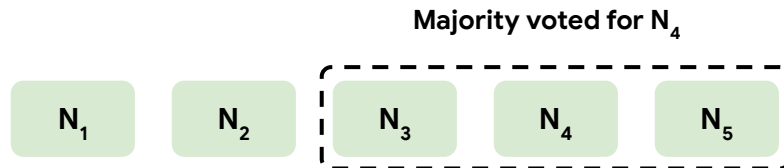
Slide is borrowed from author's presentation, restyled

# Leader election

Become candidate → Increment term, vote for self → **1** → Request votes from others

**3** → Become leader, send heartbeats

**4** → Become follower

**2** (loop back from Request votes to Increment term)

**1** After current term and vote for self persisted

**2** After timeout if cannot receive majority votes

**3** After receiving majority votes

**4** After receiving rpc from leader or with higher term

# Election correctness

**S** **Safety:** allow at most one winner per term

**L** **Liveness:** some candidate must eventually win

**Majority voted for $N_4$**

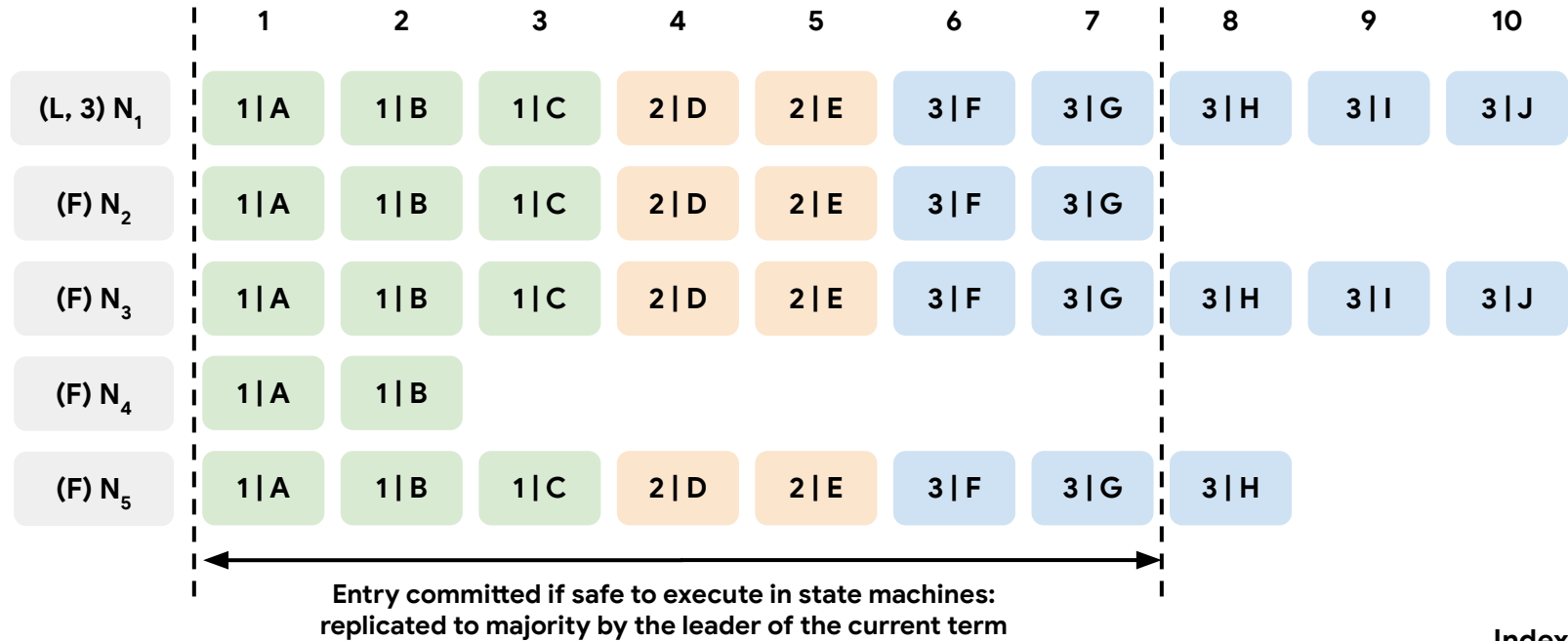$$N_1 \qquad N_2 \qquad \boxed{N_3 \qquad N_4 \qquad N_5}$$

**1** Each node gives only one vote per term (persist on disk), majority required to win election

**2** Choose election timeouts randomly in [T, 2T] (e.g. 150-300 ms), one node usually times out and wins election before others time out, works well if T >> broadcast time

# Log structure



| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|----|
| (L, 3) N$_1$ | 1 \| A | 1 \| B | 1 \| C | 2 \| D | 2 \| E | 3 \| F | 3 \| G | 3 \| H | 3 \| I | 3 \| J |
| (F) N$_2$ | 1 \| A | 1 \| B | 1 \| C | 2 \| D | 2 \| E | 3 \| F | 3 \| G | | | |
| (F) N$_3$ | 1 \| A | 1 \| B | 1 \| C | 2 \| D | 2 \| E | 3 \| F | 3 \| G | 3 \| H | 3 \| I | 3 \| J |
| (F) N$_4$ | 1 \| A | 1 \| B | | | | | | | | |
| (F) N$_5$ | 1 \| A | 1 \| B | 1 \| C | 2 \| D | 2 \| E | 3 \| F | 3 \| G | 3 \| H | | |

**Entry committed if safe to execute in state machines:
replicated to majority by the leader of the current term**

**Index**

**(Leader, Term) Node**

Slide is borrowed from author's presentation, restyled

**Term | Proposal**

# Log inconsistencies

|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| (L, 4) N₁ | 1 \| A | 1 \| B | 1 \| C | 2 \| D | 2 \| E | 3 \| F | 3 \| G | 3 \| H | | |
| (F) N₂ | 1 \| A | 1 \| B | 1 \| C | 2 \| D | 2 \| E | 3 \| F | 3 \| G | | | |
| (F) N₃ | 1 \| A | 1 \| B | 1 \| C | 2 \| D | 2 \| E | 3 \| F | 3 \| G | 3 \| H | 3 \| I | |
| (F) N₄ | 1 \| A | 1 \| B | | | | | | | | |
| (F) N₅ | 1 \| A | 1 \| B | 1 \| C | 2 \| D | 2 \| E | 2 \| P | 2 \| Q | 2 \| R | 2 \| T | |

**Raft minimizes special code for repairing inconsistencies: leader assumes
its log is correct, Normal operation will repair all inconsistencies**

**Index**

**(Leader, Term) Node**

**Term | Proposal**

# Log matching property

|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
|  | 1\|A | 1\|B | 1\|C | 2\|D | 2\|E | 3\|F | 3\|G | 3\|H | 3\|I | 3\|J |
|  | 1\|A | 1\|B | 1\|C | 2\|D | 2\|E | 3\|F | 3\|G | 4\|Q | 4\|R |  |

**1** If log entries on different nodes have same index and term: they store the same command and the logs are identical in all preceding entries

**2** If a given entry is committed, all preceding entries are also committed

**Index**

**Term | Proposal**

# Replication consistency check

|  | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| (L, 4) N₁ | 1 \| A | 1 \| B | 2 \| C | 3 \| D | |

| Before | (F) N₃ | 1 \| A | 1 \| B | 2 \| C | | |
|---|---|---|---|---|---|---|
| After | (F) N₃ | 1 \| A | 1 \| B | 2 \| C | 3 \| D | **Success!** |

**Each append from leader includes <index, term> of entry preceding new one(s), follower must contain matching entry; otherwise it rejects request leader retries with lower log index**

**Index**

**(Leader, Term) Node**

Slide is borrowed from author's presentation, restyled

**Term | Proposal**

# Replication consistency check

# Replication consistency check



Slide is borrowed from author's presentation, restyled

(Leader, Term) Node

Index

Term | Proposal

# Leader completeness

|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| (F) N$_1$ | 1\|A | 1\|B | 1\|C | 2\|D | 2\|E | 3\|F | 3\|G | 3\|H | |
| (F) N$_2$ | 1\|A | 1\|B | 1\|C | 2\|D | 2\|E | 3\|F | 3\|G | | Electing leader for term 4 |
| (F) N$_5$ | 1\|A | 1\|B | 1\|C | 2\|D | 2\|E | 2\|P | 2\|Q | 2\|R | |

**1**   **Once log entry committed, all future leaders must store that entry**

**2**   **Nodes with incomplete logs must not get elected: candidates include index and term of last log entry during vote. Voting node denies vote if its log is more up-to-date. Logs ranked by <lastTerm, lastIndex>**

**Index**

**(Leader, Term) Node**

**Term | Proposal**

# Hardened Raft

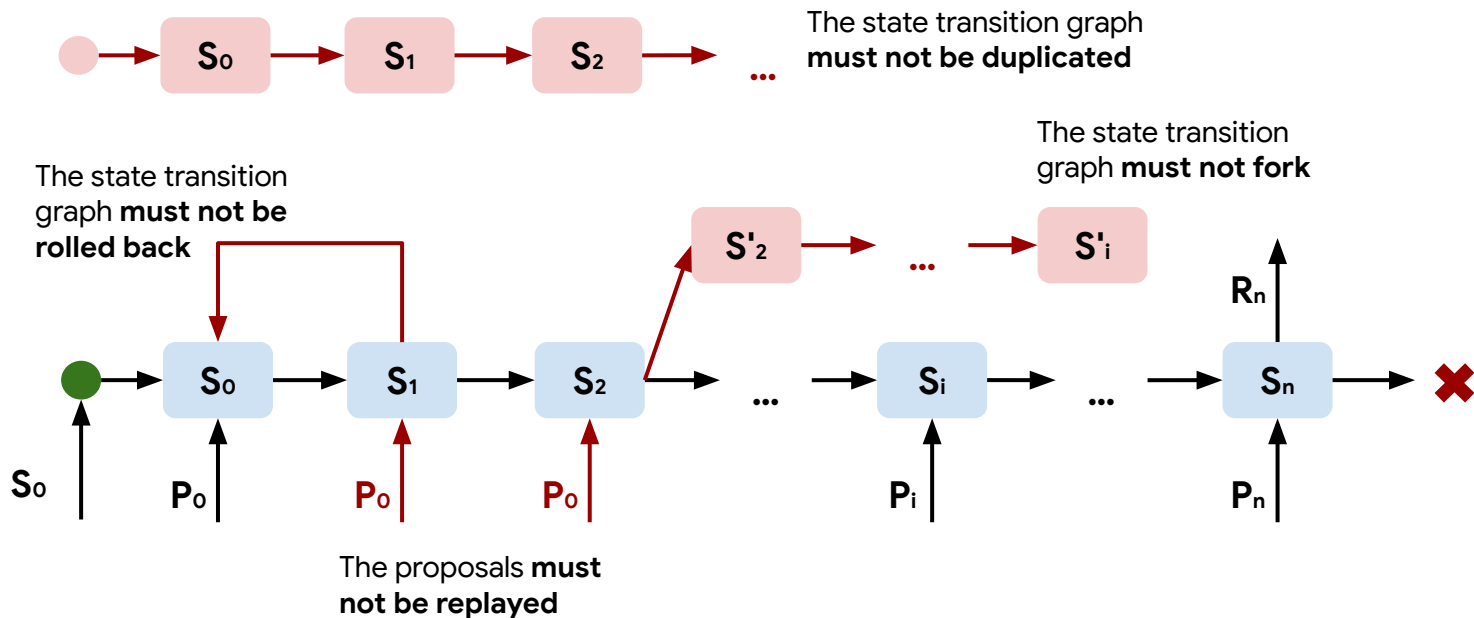Applications add app specific guarantees (e.g. user level DP for analytics)

# Protect *replicated state machines from rollback, fork or replay.

Must not forget state transitions

Must not allow identical instances to run concurrently

Must not allow inputs replay

# Attackers may control network and TEE instance execution, substitute persisted state

# Vectors of attacks

## State persistence

Raft saves its state to the persistent storage for crash recovery, including the current term, the most recent vote and log.

## Communication

Raft communicates with other nodes to maintain the consensus of the distrusted system, including leader election, replication and snapshot distribution.

# Vectors of attacks

## Node identity

Raft implementation (full details can be found in this repo) identifies nodes participating in its cluster by 64 bit integers, where node identity is used for the purposes of communication and membership. Raft relies on the fact that no two different nodes can share the same identity.

## Log integrity

Raft performs log matching during replication using log entry index and term to identify a common log prefix. Raft doesn't perform log entry contents check.

## Membership changes

Raft supports two mechanisms to change cluster membership, namely simple and joint. Both mechanisms rely on the replication mechanism to achieve consensus on the current cluster configuration.

# TEE guarantees and properties

## Tampering

TEE preserves its integrity and confidentiality during its computation, by leveraging techniques such as memory isolation, encryption and remote attestation.

## Cloning

TEE instantiates with unique encryption and signing key pairs. It is impossible to clone and run TEE instance with its encryption and signing keys.

*Machines running TEE instances are assumed to be diskless

# Raft guarantees and properties

## Safety

At most one leader can be elected in a given term (**election safety**), if a server has applied a log entry at a given index to its state machine, no other server will ever apply a different log entry for the same index (**state machine safety**).

## Communication

Raft does not require ordered or reliable message delivery for correctness (**communication safety**), allowing attacker only affect computation availability through message delay or drop.

## Log

A leader never overwrites or deletes entries in its log (**log continuity**), if two logs contain an entry with the same index and term, then the logs are identical in all entries up through the given index (**log matching**), committed entry present in logs of all leaders from higher terms (**log completeness**).

# Attacks mitigations

## Lifetime

Trusted Computation is **alive as long as the majority of the Raft cluster is alive**. If the majority ceases to exist, the Trusted Computation terminates (there are no restarts or recovery, the state is lost).
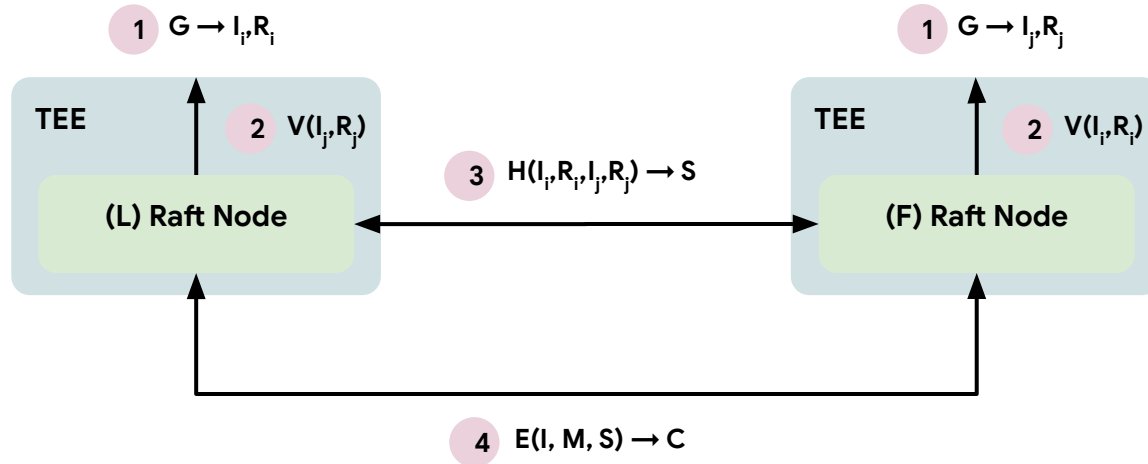
## State persistence

The machines hosting TEEs are diskless, storing in distributed persistent storage will increase cost and complexity, decrease performance during normal operation. **Raft state is not stored in the persistent storage outside of the TEE. Raft state only resides inside of the encrypted TEE memory.**

## Node Identity

Raft node identity must mean the same thing to all participants and must be resistant to impersonation. **Raft node identity is derived (hashed) from signing key pair (Sp / Sk). Raft node uses attestation report to prove it owns Sp / Sk and hence can be addressed during communication.**

# End to end communication encryption to prevent tampering and snooping



**1** $G \rightarrow I_i, R_i$

**1** $G \rightarrow I_j, R_j$

**TEE**

**2** $V(I_j, R_j)$

**TEE**

**2** $V(I_i, R_i)$

**3** $H(I_i, R_i, I_j, R_j) \rightarrow S$

**(L) Raft Node**

**(F) Raft Node**

**4** $E(I, M, S) \rightarrow C$

**1** Generate node identity and attestation

**2** Mutually verify behavior and identity

**3** Perform hybrid public key encryption handshake to agree on secret

**4** Encrypt messages with agreed secret

# Replicated log integrity checks
# to prevent log substitution attacks

**1** $G \rightarrow H_0$

$H_0$

$S_1 \leftarrow P_0(S_0)$

$H_1$

$S_2 \leftarrow P_1(S_1)$

$H_2$

$S_n \leftarrow P_n(S_{n-1})$

$H_n$

**2** $H_1 \leftarrow H(H_0, S_1)$

**2** $H_2 \leftarrow H(H_1, S_2)$

...

**2** $H_n \leftarrow H(H_{n-1}, S_n)$

**1** Generate replicated log
identity during bootstrap

**3** Maintain log integrity
through chain of hashes

Platform

Platform is suitable for customer self-hosting with straightforward dependencies

Platform enables multi-tenant setup where customer owned workers are managed by team provided coordinator

# Enable self-hosted and managed scenarios, minimal implementation surface.

Platform only requires customer to implement application actor and worker (assembled from provided components)

# Platform provides Raft nodes execution, placement and membership management, discovery and communication

# Zoom in on execution



**Consumer Endpoint**

**Node Endpoint**

4 Proposals

3 Replication, snapshot

**Executor**

2 Replication, snapshot membership, proposals, timing

**Event Loop**

**TEE**

**Raft Node**

1 Lifetime, membership

**Agent**

5 Replication, snapshot

**Comms**

1 Instructs to create and destroy nodes, maintain Raft cluster

2 Runs event loop driving execution in TEE, events are replication, snapshot, membership, proposals, timing

3 Receives replication and snapshots from leader node

4 Receives consumer proposals for execution

5 Sends replication and snapshot to follower nodes

# Zoom in on **membership management**

**Coordinator (R)**

**Worker (E)**

## Scheduler

**3** Convergence

| Nodes state | Job state |
|---|---|
| Cluster state | Placement map |

**1** Runtime state

**Agent**

**4** Lifetime, membership

**2** Desired state

**Control**

**Executor**

**Raft Node**

**1** Report nodes state and cluster membership

**2** Load desired job state

**3** Converge runtime and desired state by creating and destroying nodes, changing their membership in Raft cluster

**4** Issue node lifetime and membership commands

# Zoom in on discovery and communication

**Coordinator (R)**

**Worker (E)**

**Scheduler**

(2) **Leader resolution**

**Cluster state**

**Placement map**

(3) **Node map**

**Agent**

(1) **Node cluster view**

(4) **Node map**

(5) **Leader lookup**

**Control**

**Executor**

**Raft Node**

**Comms**

(1) Executor reports node's view of the Raft cluster

(2) Scheduler resolves which view is most up to date based on Raft term

(3) Scheduler builds and distributes node map that shows where nodes are placed, Raft cluster status

(4) Comms uses node map to enable cross node communication, clients to make proposals

(5) Lookup leader identity and location to issue proposals

# Key Value Store

# Core principles

## State and computation offload

State persistence must be offloaded to a persistent blob storage where the data is encrypted at rest (trading expensive memory for cheap persistent blob storage). State processing must be offloaded to stateless elastically scalable service (trading strict serial execution for relaxed parallel execution).

## Coarse reads and writes

The coarse reads and writes are crucial for the cost amortization. Therefore data is organized into batches that are processed as an atomic unit, thus reducing overhead associated with the number of stored blobs and the privacy state maintained by the replicated state machine.
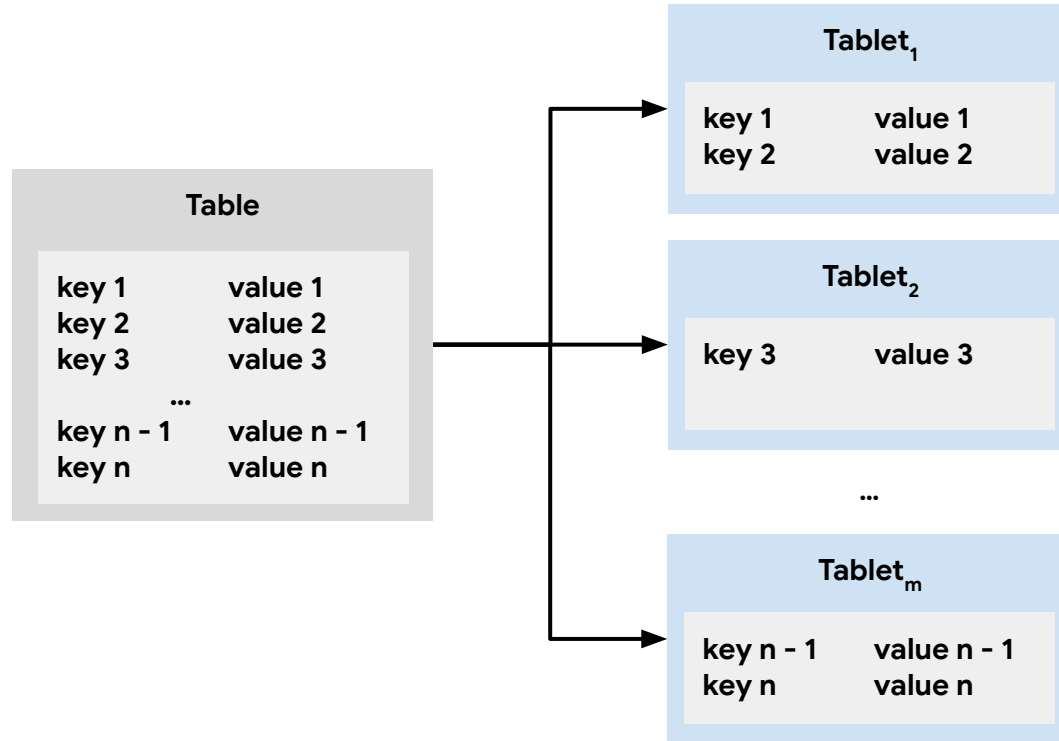
# Core principles

## Batch requests with write through cache

Requests batching serves as a bridge between fine and coarse interface while write through cache further amortizes the costs. Note that individual request latency is traded (increased) for the overall throughput (increased). Coarse operations lessen the load on the replicated state machine that are prohibitively expensive for the fine operations.
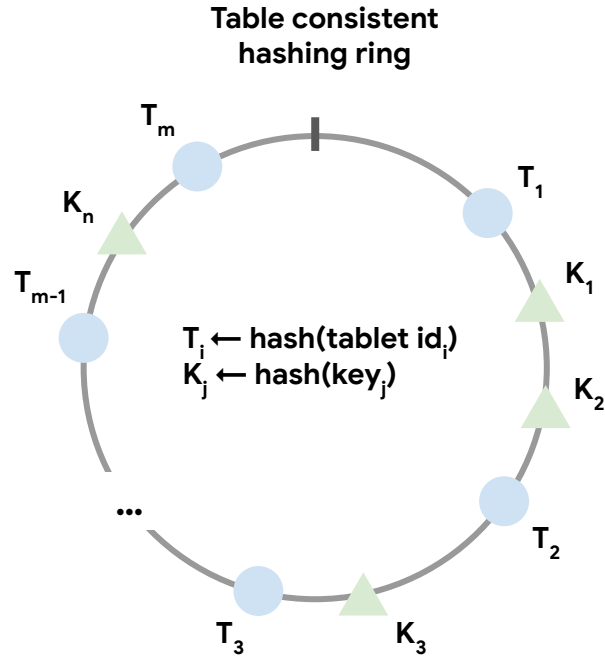
## Replicated rollback protected metadata store

Replicated state machine running inside of the TEEs provides a rollback protected replicated store that maintains metadata for the data batches. Specifically data batch is considered updated only after metadata update is committed to the metadata store.

# Split tables into smaller atomic units called tablets

**Table**

| | |
|---|---|
| **key 1** | **value 1** |
| **key 2** | **value 2** |
| **key 3** | **value 3** |
| **...** | |
| **key n - 1** | **value n - 1** |
| **key n** | **value n** |

**Tablet$_1$**

| | |
|---|---|
| **key 1** | **value 1** |
| **key 2** | **value 2** |

**Tablet$_2$**

| | |
|---|---|
| **key 3** | **value 3** |

...

**Tablet$_m$**

| | |
|---|---|
| **key n - 1** | **value n - 1** |
| **key n** | **value n** |

# Assign keys to tablets using consistent hashing

**Table consistent hashing ring**



$T_i \leftarrow hash(\text{tablet id}_i)$
$K_j \leftarrow hash(\text{key}_j)$

**Keys assigned to the closest clockwise tablet**

# **Tablet Store** manages metadata, **Tablet Data Storage** manages data

**(TEE, replicated) Tablet Store**

Tablet Metadata

Provides rollback protected replicated store for the tablet metadata

Metadata describes data in persistent storage

**(Blob Store) Tablet Data Storage**

Tablet Data

Stores actual tablet data

Load from and commit changes to the tablet metadata into rollback protected store

Cache holds recently used data and metadata to amortize costs

**(TEE, non-replicated) Tablet Cache Node**

# **Tablet Cache Node** translates key ops into tablet ops, provides write through cache

**(TEE, replicated) Tablet Store**

Load from and commit changes to the tablet metadata into rollback protected store

**(Blob Store) Tablet Data Storage**

Cache holds recently used data and metadata to amortize costs

**(TEE, non-replicated) Tablet Cache Node**

Workload issue key operations, expects cache to translate into tablet operations

**(TEE, non-replicated) Workload**

Provides write through cache for tablet data and metadata, batches key operations into tablet operations

**Tablet Data LRU Cache**

**Tablet Data**

**Tablet Metadata**

2024

Thank you

G