- Distributed storage is a key abstraction ①
- What should be the interface?
- How should it work?



why is it hard? ③

- High performance ⇒ sharding
- Many servers ⇒ constant faults
- fault tolerance ⇒ Replication
- Replication ⇒ Potential inconsistencies
- Consistency ⇒ network chit chat
  ⇒ Lower performance

what do we mean by consistency? ①

- Ideally, same behavior as a single server
- writes one at a time (even if concurrent)
- reads latest write.

* Don't want to talk about internals of storage

* Judge by behavior seen by clients.

Create an illusion of single node storage.

---

Example    ②

C1:    ├─ Wx10 ─┤

C2:    ├─ Wx20 ─┤

                        ├─ Rx? ─┤

C3:

                        ├─ Rx? ─┤

C4:

- Either 10 or 20, but same.

- Single server has poor FT/scalability

---

Bad replication design    ①



WX10        Rx10

q ──────→ S1 ←--- C3
  \  10x10 ↗    ✗
   ✗
  WX20  ✗
         Rx20
C2 ──────→ S2 ←--- C4
    WX20

• C1 crashes before sending to S2

• S1 crashes after C3 reads, before C4 reads.

Suddenly x=20.

Consistency problems

$C_1$ ⊢$wx10$⊣

$C_2$ ⊢$w_x 20$⊣

⊢$R_x 10$⊣ $C_3$

⊢$R_x 20$⊣ $C_4$

$$C_1 \rightarrow S_1 \rightarrow C_3$$
$$C_2 \nearrow\!\!\!\!\times S_2 \rightarrow C_4$$

- CR Goals -

        ~~FT~~ ~~Distributed~~ Key value store

- Data fits in memory / disk
- Many concurrent clients

- Consistent view of storage.

        we want that

→         all execution histories

are <u>linearizable</u>.

· One can find a total order of all ops

    * matches real-time for non-overlapping
        ops

    * read sees last write

(4)

C1: |—•—Wx1—————| |—•——Wx2——|

C2: |—•—Rx1——|          actual ops happens
                        anywhere in
C3:        |—Rx2——•|         between

                        client recieved
                        response

client    start
sent      time
request

⇒) linearizable

---

happens before in real time          (4)

C1    |—Wx1—| |—Wx2—|

C2:        |—Rx2—| G      happens before in
                          read-write
c3:            •|—Rx1—|

cycle ⇒) NOT linearizable

---

(4)

|—Wx0—•| |—Wx1—•—|

                |—•—Wx2—|

              |—Rx2—•—||—Rx1—•—|

⇒ linearizable.

In concurrent writes, write that started
later can be executed first.

- Linearizability is a **safety** property ← Bad things never happen

C1: ⊢— Rx —— ...
C2: ⊢— Wx —————— ...

} Trivially satisfied by a storage system that does nothing!

⇒ we also want **liveness** (even w/ fault)
  ↳ Good things eventually happen

---

② (2)

- **Liveness**

  - Good things eventually happen.

Research. safety-
  * will not publish incorrect results/ low quality research

Liveness
  * will eventually publish

---

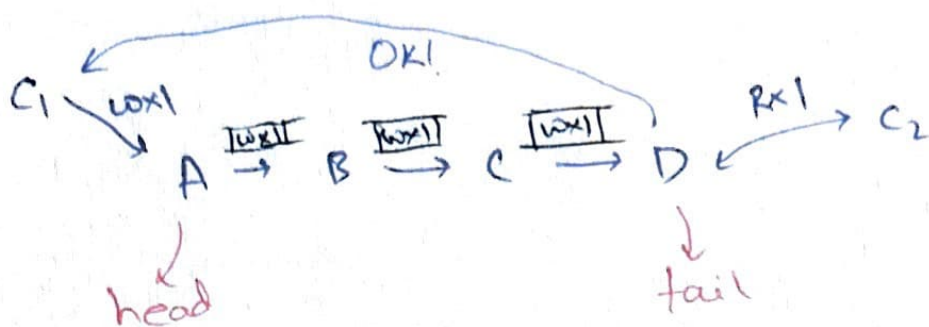**Safety, liveness in concurrency.** ②

**Safety:**

- Mutual exclusion

- No deadlocks

**Liveness**

- Each process will eventually get the shared resource (starvation freedom)

- Chain Replication



clearly linearizable: in terms of tail
server

- Why TLA+ (Temporal logic of actions)

  - As storage systems grow in complexity, reasoning that they behave like a single system becomes very difficult

  - Adding optimizations to existing systems often can break guarantees.

- Days of debugging can save hours of writing specifications.

- TLA+ lets us formally prove than one system implements another. Eg. CR⊘ implements single machine storage (linearizable)

- Prove safety / liveness

---

BIG IDEA:

TLA+ helps convert implementation into mathematical implication

$$CR \Rightarrow SS$$

---

* First 1/2ʳᵈ of course
- Learn TLA+. Prove that CR implements SS.
- Very beautiful concepts - stuttering, refinement, actions, ...
* Second 1/2
- Discuss distributed storage systems.
Groups of 2.