- If single computer is enough dont use DS.
- Motives for DS:
① Parallelism
② Fault tolerance : If one fails, another is backup
③ Physical reasons : systems inherently are physically distributed
④ Security : require isolation

- Unexpected failure patterns ⎫ Basic
- Partial failures              ⎬ Challenges
- Performance

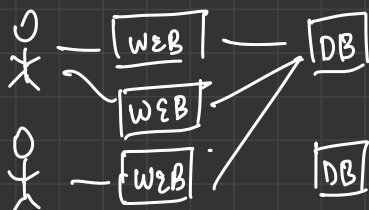- Infrastructure
- Storage                    ⎫
- Communication   ⎬ Abstractions ─┐  → Implementations:
- Computation           ⎭                    - RPC, threads - way of
                                                                            structuring
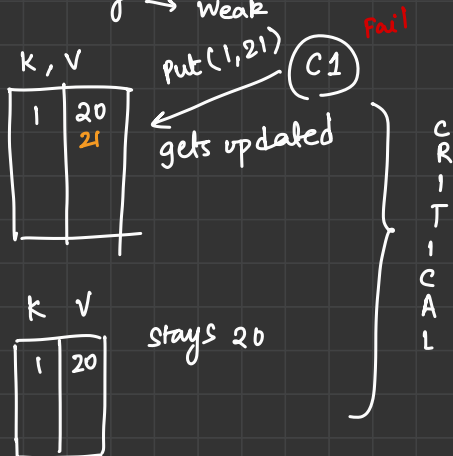                                      - concurrency        concurr. ops
                                          control

- Performance
- High level goal : scalable speed up    2x machines =
                                                                    2x throughput

- Fault Tolerance  ] Big tools
  - Availability | clever ways of avoiding to write non-volatile storage
  - Recoverability | management of replication

- Consistency → strong: guarantee of most recent write is expensive
              → weak

K , V    put(1,21)  (C1) Fail

| 1 | 20 |
|   | 21 |

gets updated

$\left.\begin{array}{c}C\\R\\I\\T\\I\\C\\A\\L\end{array}\right.$

- Replication for fault tolerance should have independent failure probabilities.
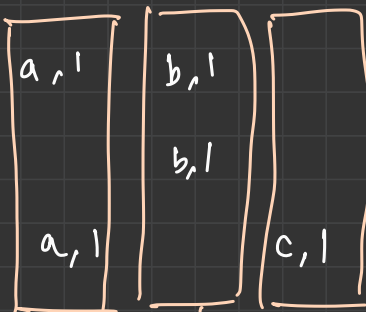
  ↓

  might make communication expensive

K   V    stays 20

| 1 | 20 |

Jan 22, 2023

→ map Reduce

Input File 1 → map
Input File 2 → map
Input File 3 → map

| a , 1 |   | b , 1 |
|       |   | b , 1 |
| a , 1 |   | c , 1 |

map (k, v)
split v into words
for each word w
  emit (w, 1)

reduce (k, v)
  emit ( len (v))
           ↑
         array

  ↳ Reduce → (c, 1)
  → Reduce → (b, 2)
  → Reduce → (a, 2)