

• DS = collection of computers that com via networks

↳ Inherent distributed

- ↳ necessary
- ↳ sufficient

↳ Better reliability

- ↳ one node fails \Rightarrow system can't work

↳ Better performance (fast)

↳ Solve bigger & yrs problems

- ↳ single node cannot handle all data / processing capacity

↳ efficient solve

Problems \rightarrow coms Fail

- \rightarrow node crash
- \rightarrow Faults: random / coord
- \rightarrow security & privacy

Properties

- Architecture
- Process
- Coms
- Coord
- Naming
- Consistency & Replicat
- Fault Tolerance
- Security

Design goals

1) Resource sharing

- achieve effectiveness
- manage costs

2) Distribut - transparency

access / locat / relocat / migrat / replicat / concurrency / failure

3) Openness

components easily integrated onto other systems

4) Dependability

- ↳ Availability: readiness usage
- ↳ Reliability: continuing service
- ↳ Safety: low catastrophe prob
- ↳ Maintainability: ease repair

5) Security

confidentiality / integrity / trust

6) Scalability

↳ Size scalability

= n° users / process / nodes

- limited comp / storage / network capacity
- scale up vs scale out

↳ Geographical scalability \rightarrow latency: Timeouts (1)

- hard LAN to WAN

- high latency some apps / protocols may not work

- unreliable WAN links

↳ Admin

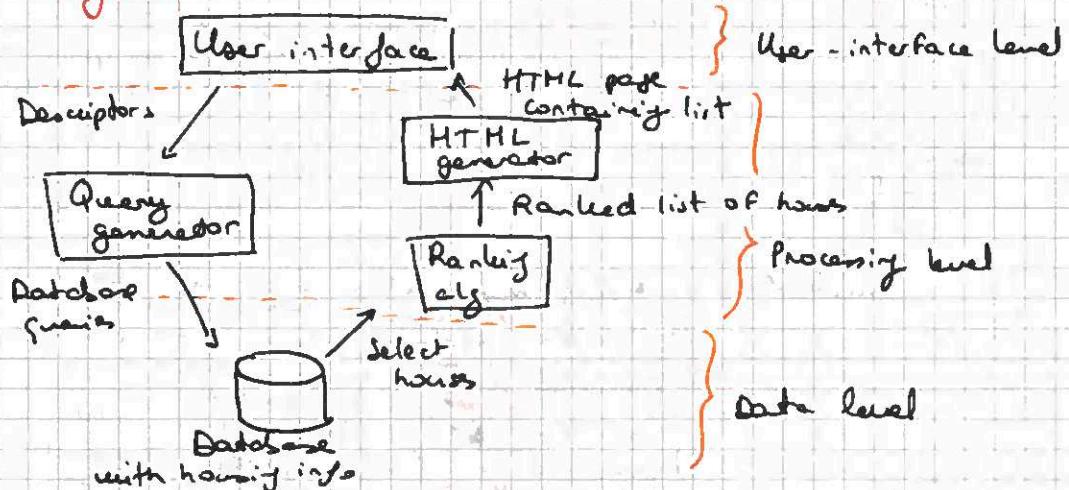
- Conflicting policy concerning usage
- Many systems avoid this.)

Architecture

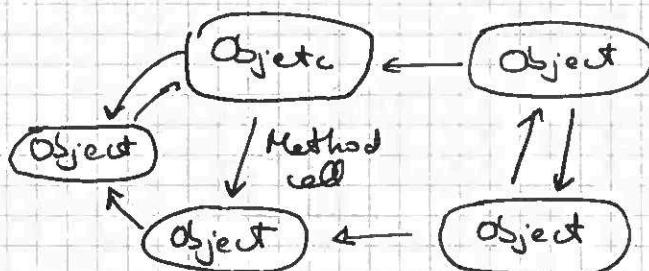
Style

- components with well-def interface
- how components are connected
- data exchanged b/w components
- how components & connectors are jointly configured
= mechanism mediates comp, coord, coop among components

Layered

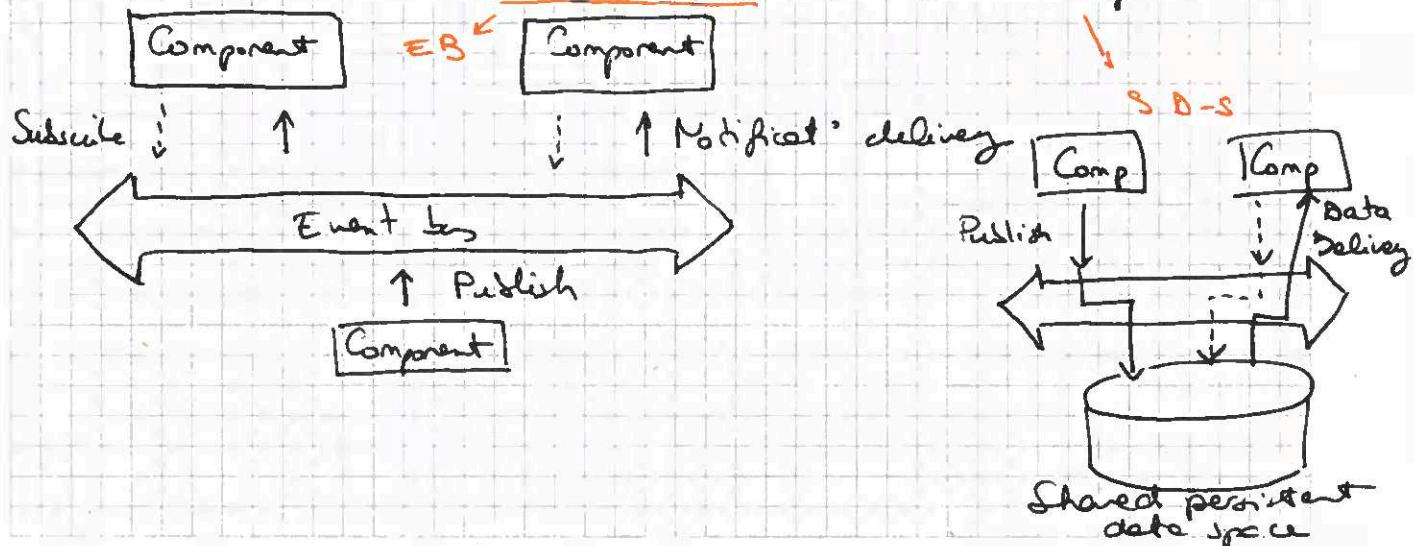


Service oriented / object oriented, microservices



Publish - Subscribe

System = collect + autonomous operating process
sharing report, processing & coord
Event-based / shared data-space

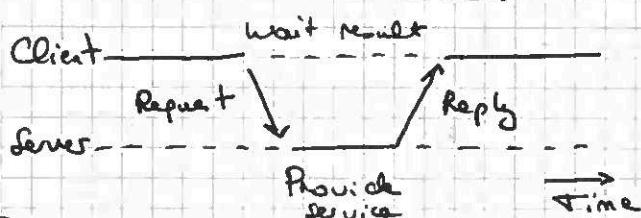


Systems

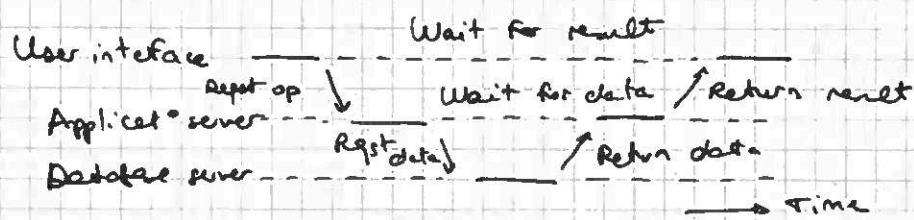
→ Centralized

Client - Server
process offering services

2tier

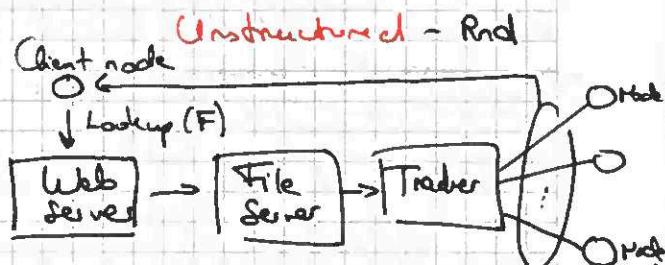
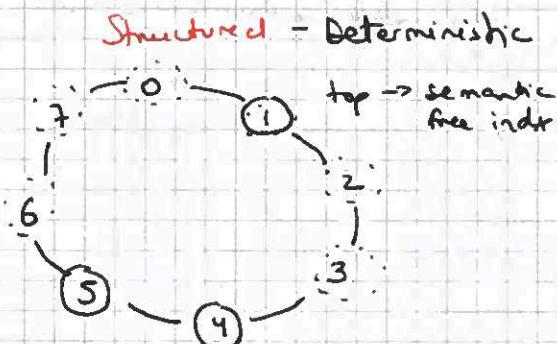


3tier



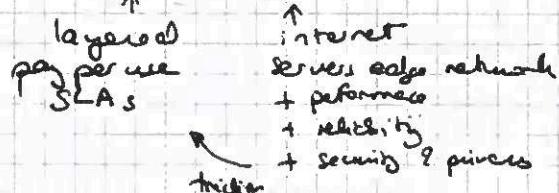
→ Decentralized Peer-to-peer

Client server physically split into equiv parts which operate on own share data + load balancing



→ Hybrid

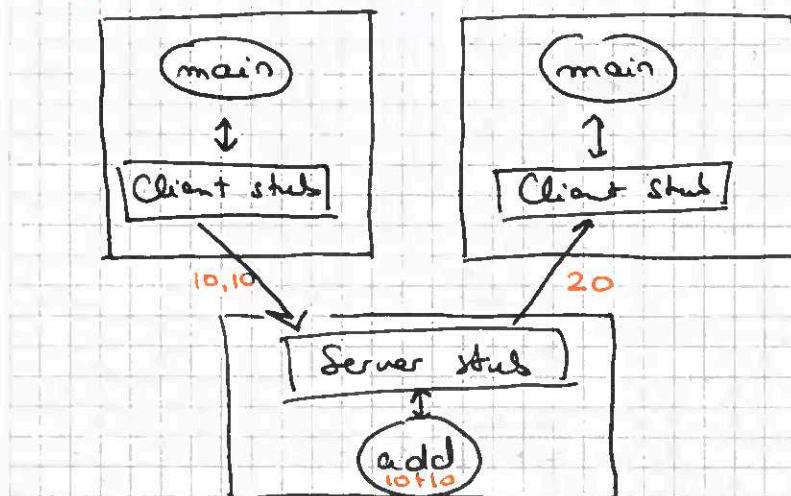
ex: Cloud, Edge, Blockchain



- * RPC => allow remote service to be called as procedures
 - transparent w.r.t
 - ↳ locat
 - ↳ implementat
 - ↳ lang

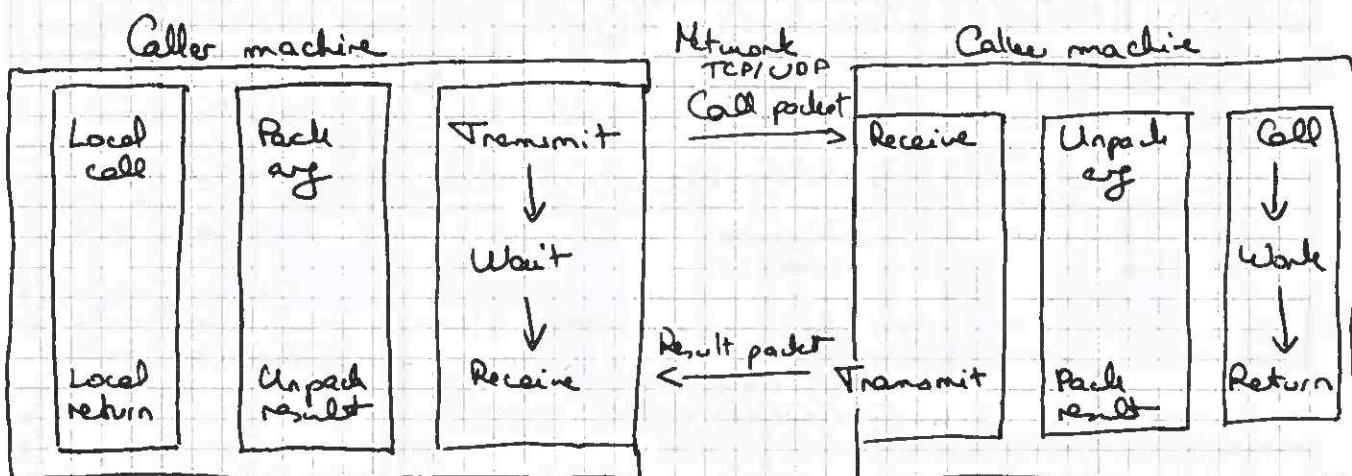
Goal: Make dist comp look like centralized comp

≈ Programs can call procedures on other machines



Challenges

- separate callee & caller address space
- machines may be different
- thousands of procedures
- client & server may fail (indep)



How do we handle scaling?

- Hide coms latency < asynchronous coms (post-reply behavior)
 - ↳ delegate work
- Partitioning & Distribute = split data & spread - decentralized - cache data
- Replicat': drawback consistency

Dependability requirement \leftarrow Availability
 \leftarrow Reliability : continuity service
 \leftarrow Safety
 \leftarrow Maintainability

- * MTTF = Mean time to failure ↑
- * MTTR = ... repair ↓
- * MTBF = Mean time b/w failures ↑

$$\text{Availability} = \frac{\text{MTTF}}{\text{MTBF}} = \frac{\text{MTTF}}{\text{MTTF} + \text{MTTR}}$$

- * Failure = syst where component(s) not working
- * Error = part component that can lead to failure
- * Fault = cause of error

- * Fault tolerance = syst continues to work under **Faults**

transient

 intermittent
 permanent

Two general's problem "no solut"

only able to com with each other by sending messages through enemy territory

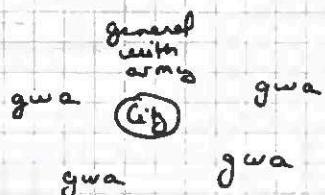
- : no guarantee delivery
- : no acknowledgement chain

- probabilistic agreement
- retry mechanism
- multiple channels

Bizantine general's problem

$3F + 1$ nodes tolerate F malicious nodes

$$\frac{3F+1}{\text{total}} - \frac{F}{\text{bad}} = \frac{2F+1}{\text{majority}} \hookrightarrow \text{good}$$



friends must agree pbn
don't know bad

bad may collude

Network Behaviour

Arbitrary: malicious adversary may eavesdrop, modify, ...

↓
Fair loss: message may be lost, duplicated, re-ordered, ...

↓
Reliable: message received if sent, m can be re-ordered.

Halting Fail

- stop: crash failure, reliably detects failure
- noisy: eventually reliably detect failure, noisy behavior] fail exp
- silent: omis or crash failure (client won't know what's wrong)
- doze: benign failures (no harm)
- arbitrary (byzantine): with malicious failures

multiple
clients

failure

Timing constraint

2 delivery times

→ Synchronous: process execute speeds bounded
can detect omis & timing failures

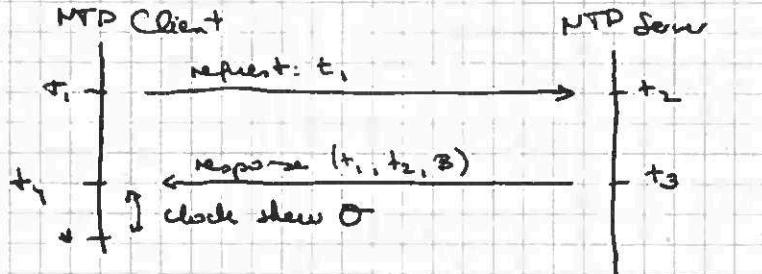
→ Asynchronous: no time constraint
cannot reliably detect failures

→ Partially synchronous:

Clock (computers) = quartz oscillators \rightarrow drift

* Clock skew: difference b/w 2 comp clocks (asymptotic)

\Rightarrow Synchronisation - External: external server bounded by D : $|C(i) - s| \leq D$
 Internal: $|C(i) - C(j)| \leq D \Rightarrow$ a bound error



$$\text{Round trip time } RTT = (t_3 - t_1) + (t_3 - t_2)$$

$$\theta = t_3 + \frac{RTT}{2} - t_4$$

Synch issues

Time relativity

Variable com latency \Rightarrow Absolute time NOT A GOOD IDEA

Consistency vs time

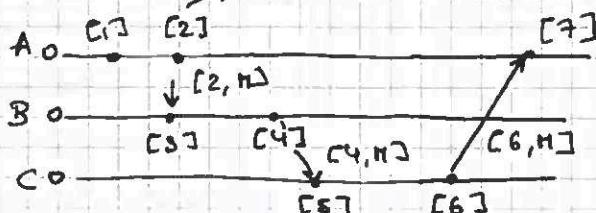
Faults

Logical Clocks order events (consistency)

- Rules
 - \rightarrow each process maintains counter val
 - \rightarrow each process increments counter when send / execute
 - \rightarrow send message event counted (timestamp)
 - ~~\rightarrow receive message event +1 local counter~~

$$\max(\text{local counter}, \text{message counter (timestamp)}) + 1$$

- $a \rightarrow b \Rightarrow \text{val counter}(a) < \text{val counter}(b)$ \Rightarrow Lamport Clocks do not guarantee ordered or unifl concurrent events



\rightarrow L_clock = single scalar val to track causal ordering

Vector clocks

maintain n vals for n nodes in syst. Vector stamp of an event

$$v(a) = (t_1, t_2, \dots, t_n)$$

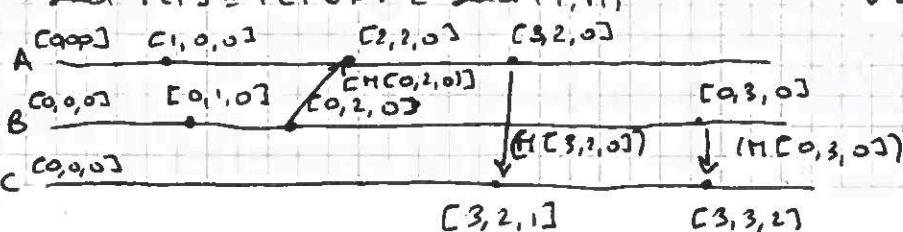
t_i : n° events on node i :

each node has current timestamp T

Rules

- \rightarrow initializat? $T = (0, 0, \dots)$
- \rightarrow event occurs $T[i] = T[i] + 1$ ~~& send (T, M)~~
- \rightarrow send $T[i] = T[i] + 1 \in \text{send}(T, M)$

$$\begin{aligned} \rightarrow \text{receive } TC[j] &= TC[j] + 1 \\ TC[k] &= \max(T[k], T'[k]) \\ \forall k \neq j, k \in C_1, \dots, n \end{aligned}$$



→ Global snapshot = Global state, comprises → indiv. state of each process
 → indiv. state each uses channel

Requires

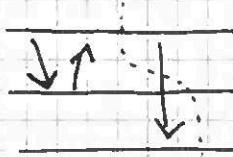
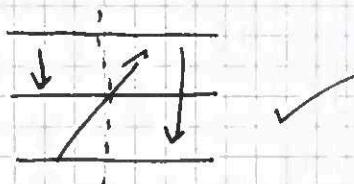
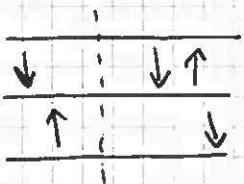
- ↳ each process must record self
- ↳ record content coms channels
- ↳ no pause system or interference with normal func'
- ↳ collect in distributed manner

- Support checkpoints
- garbage collect
- deadlock detect
- transient computes

Consistent cut

receiver include :: so should at sender

events not so



X : events not causally related
 (consistent events should be in same cut)

thinkLab L5: Consensus with Paxos

IBM

Paxos agrees each client q

Propose ① Propose val

Prepare

② Acceptor: ok as long as $n > n_h$

Accept ③ Proposer got majority \rightarrow here is the val

④ Reply acceptors to proposer ✓, otherwise try again (terminate)

Proposer: proposed n

Acceptor: if $n > n_h$

$n_h = n$
if no prior proposal accepted

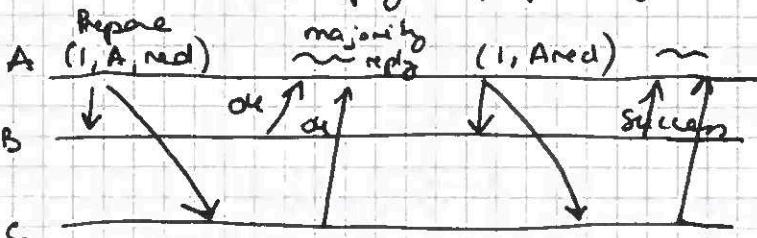
reply $\langle \text{promise}, n, \text{NULL} \rangle$

else

reply $\langle \text{promise}, n, (n_a, v_a) \rangle$

else

reply $\langle \text{prepare failed} \rangle$



Priority safety \rightarrow liveness: not guaranteed, e.g.: competing proposers
only one val chosen \rightarrow system will eventually make progress
val remains consistent

Raft alt to Paxos, design simpler: optimized for (log) append

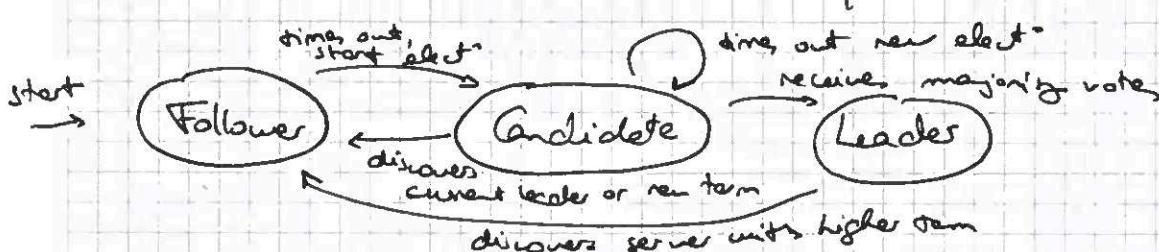
agree on new leader

Roles

- Follower - passive: respond reqts
- Candidate - used to elect new leader
- Leader - handles all client requests

• Leader by a truth
L keeps track of entries
L updates F

L clears up
L minimizes entries
L starts
L vote uncommitted
L repeated
L ...



- terms act as logical clock
- each node stores current n = increases monotonically
- current terms exchanged while normal ops
- reqst with stale n \Rightarrow rejected

L6: Mutual Exclusiveness & Concurrency

non-dict behavior

* **Critical sect** = uninterruptable group of code (ensures atomicity)

* **Race conditions** = error dep on timing

Mutual exclusiveness → Critical sect executes uninterrupted

Deadlock condit → eliminate 1 ⇒ fix

ME HW NP CW

- Mutual exclusiveness: at least one resource in "non-shareable" mode
- Hold & wait: process holding at least one resource waiting to acquire additional resources
- No preempt: resource cannot be forcibly removed from a process holding them
- Circular wait: circular chain of processes, where each process holds a resource that the next process in chain is waiting for

Consistency ⇒ strong operations on shared data are sync without sync ops

Causal

concurrent writes may
be seen diff order on diff
machines

FIFO

writes performed by a single process are
seen by all other processes in order
they were issued

P₁ W(x)a

R(x)a W(x)b

P₂ R(x)b A(x)a] violates

W(x)a

A(x)a W(x)b W(x)c

P₃ R(x)a R(x)b

R(x)b R(x)a R(x)c

P₄

R(x)a R(x)b R(x)c ✓

wouldn't work if R(x)c before R(x)b

P₁ W(x)a

• **STRICT** absolute time ordering of all shared access
strong matters

P₂ W(x)b

• **SEQUENTIAL** all processes see shared access in same
order

P₃ f(x)b R(x)a

• **LINEARIZABILITY** sequential consistency + ops ordered
according to global time

P₄ R(x)a R(x)b

• **CAUSAL** all processes see causally related
shared access in same order

• **FIFO** all processes see writes from each other in
order they were issued

→ **Weak** = shared data can be consistent only after sync has done (release/abort)

Weak model = more scalable

L7: Replicat^o

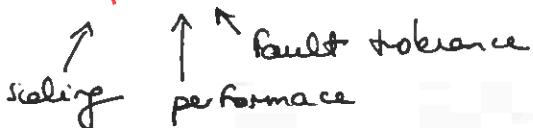
~~transactional properties~~ Atomicity - all or nothing (commit/abort)

Consistency - follow rules

Isolation - concurrent transacts don't interfere with each other

D Persistence - remember everything (commit → permanent change)

Replicat^o ← availability

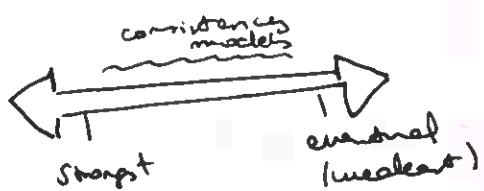


Consistency Models

Data Centric defined consistency experienced by all clients system wide

L8: Client centric consistency models

Client Centric defined consistency of data stored only from one client's perspective (diff clients may see diff data)



6 Guarantees consistency models

- 1) Read my writes: client will always see effect of its own writes
- 2) Monotonic reads: no going back in time
- 3) Bounded staleness: each guaranteed to return data to client in predefined window (temporal consistency)
- 4) Consistent prefix: client will see updates in correct order
- 5) Eventual consistency: no updates → all replicates will eventually converge
- 6) Strong consistency: all clients see same data simultaneously after write complete

Strong - see all prev writes
2-5

Trade-offs

	Consistency	Performance	Availability
S	Excellent	Poor	Poor
E	Poor	Excellent	Excellent
CP	Okay	Good	Excellent
BS	Good	Okay	Poor
MR	Okay	Good	Good
RMW	Okay	Okay	Okay

Eventual - see subset prev writes

0-0, 0-1, 0-2, 0-3, 0-4, 0-5,
1-0, 1-1, 1-2, 1-3, 1-4, 1-5,
2-0, 2-1, 2-2, 2-3, 2-4, 2-5

Consistent prefix - see initial seq of writes

0-0, 0-1, 1-1, 1-2, 1-3, 1-4, 2-4, 2-5

Bounded staleness - at most one ring out of date
→ how out of sync U can tolerate

2-3, 2-4, 2-5

Monotonic reads - one client reads w/ all subsequent reads will reflect same or more recent vals, no "time travel" → if we read something older we jump back

1-3, 1-3, 1-4, 1-5, 2-3, 2-4, 2-5, _{jump back}

Read my writes - client will always see effects of its own writes & reads

writes: 1-5

reads: 0-0, 0-1, 0-2, 0-3, 0-4, 0-5, 1-0, 1-1, 1-2, 1-3, 1-4, 1-5, 2-0, 2-1, 2-2,
2-3, 2-4, 2-5 - every opt

L9: Network FS

* **File** = array persistent bytes that can be read/written

* **File system** = • collects files

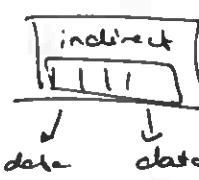
• part OS manages files

Name → access → **Inode** = file that stores data & metadata of a file
(in memory rep of file)
↳ inode n = (unif. id)
↳ path
↳ file descriptor

Multilevel inode



↓
data
↓
data



↓
data
↓
data

* **Directory** = specific file that contains other files

- stores entries in data blocks
- distinguished from other files via inode

↓
data
↓
data

* **Superblock** = starting point that stores important info about file syst (usually 1st block)

↳ n° blocks
↳ n° inodes
↳ block size

File system Types

Local

Network

NFS [pull]

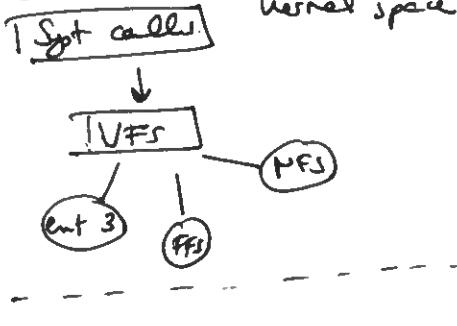
process on same machine as
access shared files

Virtual FS
virtual abstract →
superblock
inodes
file

Apps

user space

kernel space



States, protocol
+
File handle
+
Client logic

Client → file handle (returned)
↑ current offset

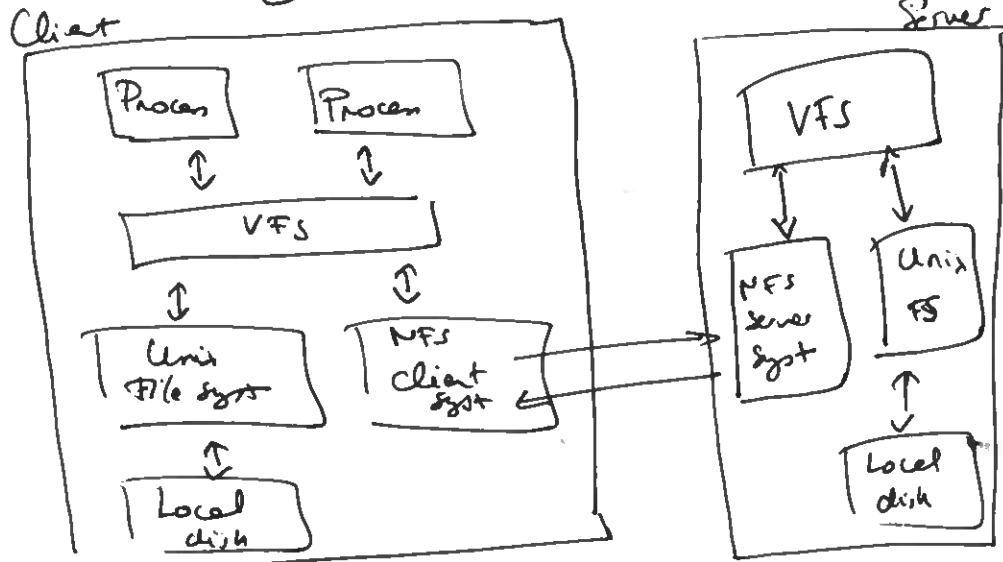
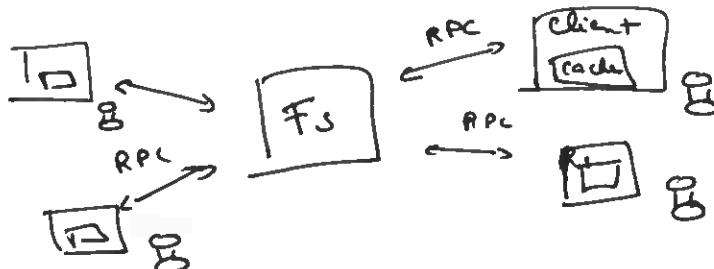
process on different machines access stored files

- fast & simple network recovery
- transparent access = normal UNIX semantics
- reasonable performance (scalable)

Network FS → protocol: stateless

(NFSv2) every reqst from client completely describes
required ops

- no client state maintained
- server can keep state for performance but not correctness
- crash/reboot → no correctness problem, just slower performance



* **Idempotent op** = same result regardless no. times called

(e.g. read, pwrite)

 ↑ write to specific offset of file (pos=*)

Non-idempot. (e.g. append, mkdir, rmdir, rm, etc.)

ff
p-point
mk - make
rm - remove
buffer - temp mem
storage
cache - comp in media
mem

* **Caching** = temp storage for freq accessed data in fast access storage

NFS problems

- **Server Caches**

↳ buffer write to improve performance

↳ write might be acknowledged before push

→ socket → don't use server write buffer: persist data to disk before write
 → use persistent memory → more expensive

- **Caching** → update visibility

... server doesn't have last ver.

... clients may see old ver.

also NFS

→ flush on close (or other time) → problem
 flushes not atomic (one block at a time)
 two client flush at once ⇒ mixed data

- **State Caching**: clients still have stale data they shouldn't read

→ clients recheck in cached copy up to date = STAT request

STAT gets last modify time stamp +
↳ problem

overloaded server repst: 90%

↳ set to every 3s

push

NFS

Goals = more reasonable semantics for concurrent file access
 improved scalability
 willing to sacrifice staleness

- whole file cached (instead of byte)

- last writer wins

State cache → Full state protocol

server tell client when data overwritten

server remembers which clients have file open right now

client verifies checksum when open file (not every read)

Client push ⇒ enrich everything cache + recheck entries

 → avoid re-reading old data

Problems

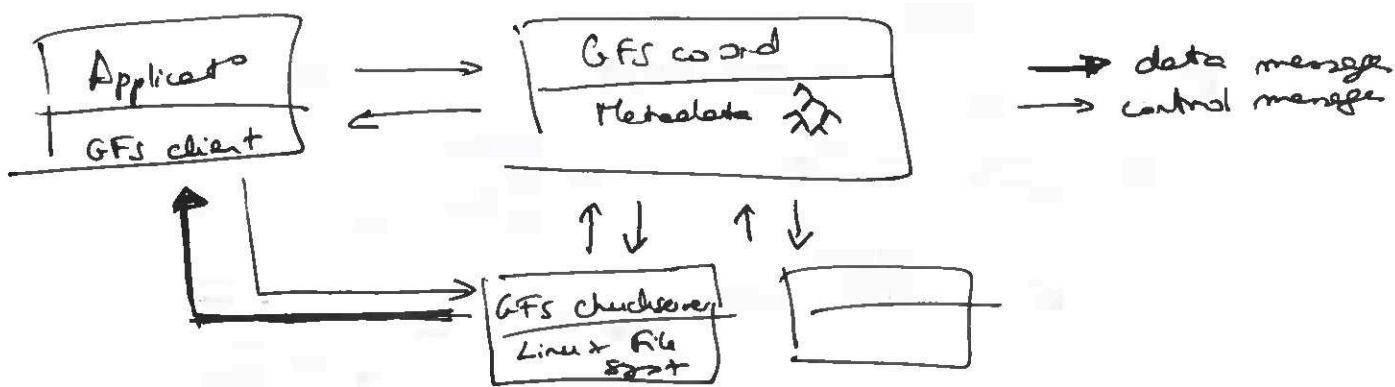
→ server → tells client recheck data (it loses track who open file)

L10 : GFS

GFS Scalable data intensive applicat'

- Goals
 - large storage
 - scalable access
 - fault tolerance
 - transparency - local
 - fault

- Assumptions
 - failure common
 - file large
 - most files appended, not overwritten
 - and writes never done (almost)
 - file access
 - concurrent access
 - needs
 - possible
 - large
 - small
 - streaming
 - files mostly read (synchronous)



Single coordinator: state replicated on backups

- ↳ Holds all **metadata**
 - name space
 - access control
 - replica location
 - check leases
 - chunk locators
- ↳ **Manages**
 - garbage collect or optional chunk
 - chunk migration
- ↳ **Fault Tolerance**: periodic com with chunks via ↗
up to 3 replicated multiple machines

Chunkserver

stores chunks on local disk

chunk size = 64 MB

32 bit checksum - detect corrupt

each chunk 3 replicas

chunk handle

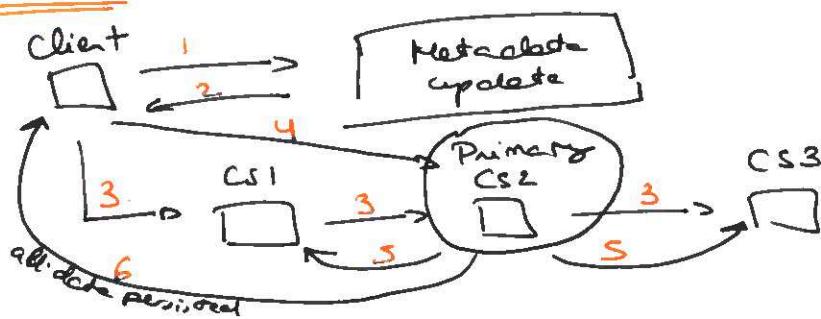
↳ global unique 64bit ↗

↳ assigned coord

↳ read / write req specify handle

Apps must cope
with stale
incorrupting
(failure)

Client write



Metadata consistency

- namespace changes are atomic (all or nothing)
- coordinator uses op log to store critical metadata changes
- log defines timeline → def order concurrent ops
- log stored on coordinator local disk & replicated
- coord recycles file system by replacing op log
- coord only replies to client after log entries safe on local disk & replicas

Data consistency

	Write	Read
Serial success	defined	defined
Concurrent success	consistent, unclef	int
Failure	inconsistent	inconsistent

def: primary tell client write succeed & no other client wrote same part

consist: all concurrent writes succeed, readers see same content, mix

inconsistent: primary doesn't tell client op succeed
→ diff readers see diff content or none

(*) secondary ^{none} do not update

Apps

- must deal with inconsistency
 - cannot tolerate duplicate → unif identif used
 - rely on checksum
- ↑ coordinator detects & removes

Handling Faults

Secondary

- primary retry
- client retry
- coord may remove cc from list (checkin/checkout)

Primary

- coord remove cc from list
- coord grant lease to any secondary

Coordinator

- replay op log → rebuid state
→ remove ops
- ask cc what they store
- wait for one lease before granting to any secondary