

disorderly-labs.github.io

Distributed Systems → fundamentally characterized by uncertainty.

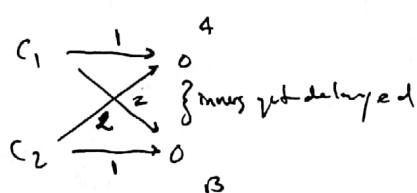
DEFINITION: DISTRIBUTED SYSTEM:

- * A bunch of computers Asynchronously communication
 - * Many computers appearing as one. → cannot make assumptions on upper limit of message time.
 - really hard to detect failed connection on an asynchronously distributed system.
 - * Parts are communicating over an inconsistent network.
- WHY
- * Reduce latency (due to locality) - replication redundancy
 - * Scalability
 - * Alternate routes - multiple choice same task redundancy

DEF: ARE non-deterministic. uncertain.

Uncertainty

1. Concurrent / Asynchrony
↳ Message ordering.

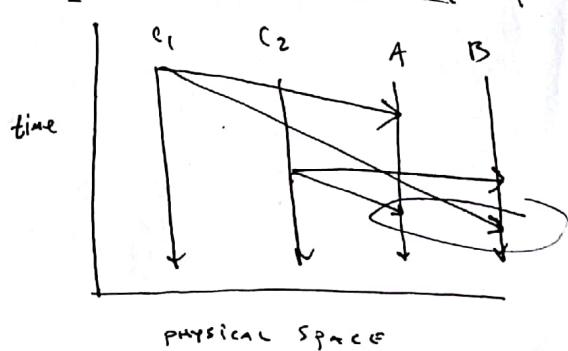


If B is a replica of A.
And some delay occurs,
A & B will have input
different messages.

non-deterministic delivery order.

2. FAILURE (Partial Failure)
↳ Components can fail, without entire system failing.

SPACE TIME DIAGRAM (LAMPORT)



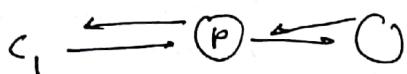
THEME 2 : time & Asynchrony.

no total Order in DS.

Partial Orders

reasoning about time & Partial Orders.

THEME 2 : FAULT TOLERANCE



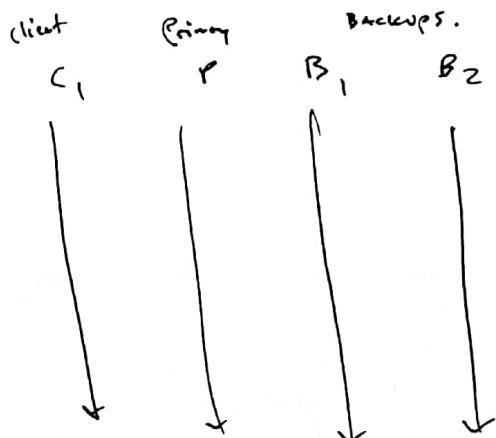
Primary replicates & waits until they signal that the stored data to respond to client.

Big concern is with primary failing.
client thinks success, no data stored.

THEME 4: Scalability & Parallelism.

THEME 3 : Consistency.

So if I can trick you into thinking I'm not a collection of computers, I'm strongly consistent.



(1)

TIME, CAUSALITY, "HAPPENS BEFORE" NOTION, LOGICAL CLOCKS, PARTIAL ORDERS.

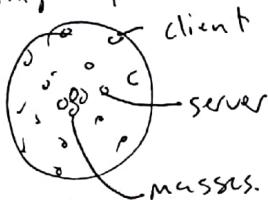
Come up w/ questions about RESEARCH PAPER.

(1) (2)

WANT DS?

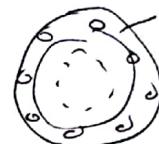
- 1. Scale-up / Scale-out } most DS want BOTH.
- 2. FAULT-TOLERANCE }
- 3. REDUCED LATENCY → LOCALITY ... }

Extrinsically Distributed Systems.
 • Nothing about algorithm is Distributed.
 • Mass is in center, pretending to be a single system.



• correctness = transparent distribution.

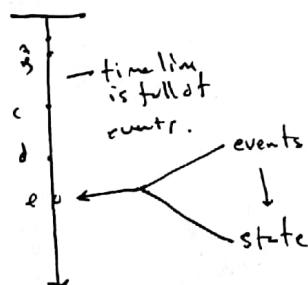
- LESS COMMON, MORE CUTTING EDGE
- INTRINSICALLY DISTRIBUTED.
- P2P systems.
- Mass is on edge
- IoT.



need to invert models of correctness.

LAMPORT Diagrams

TIME:



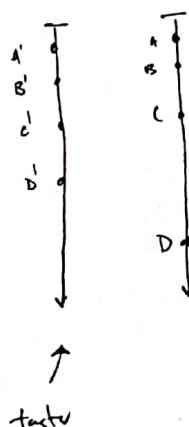
e is the set
 $(A \in B \in C \in D)$

• state is a ~~history~~^{lossy}
 representation by
 a chain of events.

usually THE CASE, given a history (set of events) we can figure out state.

(a particular time is the same as history of previously occurred events.
 → event go state → event.

causal ordering - accounting for time -



- can't distinguish lines, due to fact that the same events occurred on each line.
- consequence of Asynchronous Network Model.

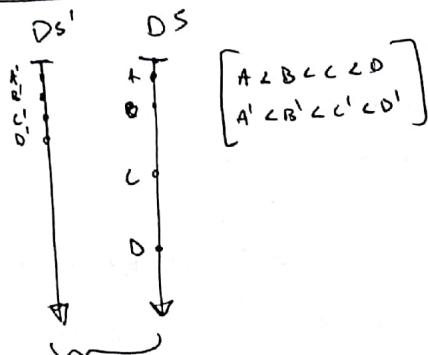
NETWORK MODELS

Synchronous

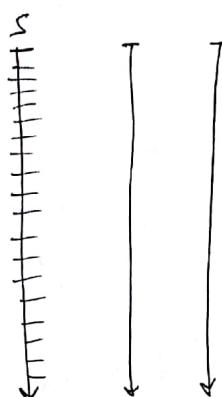
- upper bound on the latency of computation & communication.
- can only take so long: for message delivery
task completion ...

Asynchronous

- no bounds



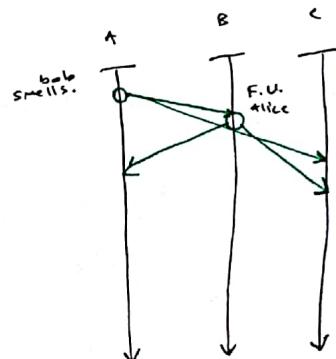
Same results.
DS' gets job
Done faster.



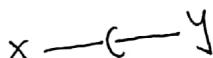
[what is order for δ ?]
[why do we need it?]

consistency

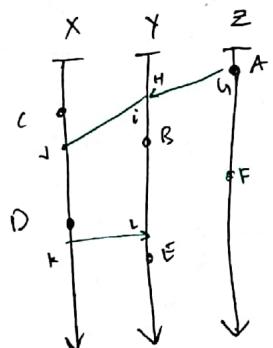
- exhibits behaviors that look the same way.
- implies something about anomalies that do not occur.



CAUSALITY



Potential Causality

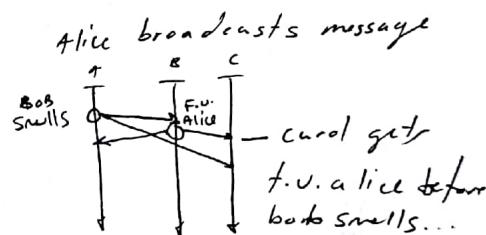


"HAPPENS BEFORE" RELATION

- Event B caused A?
- B could've caused E
- A can't cause an event
→ nothing is reachable from A.

E is set of events $\{A \prec B \prec C \prec D\}$
 $\rightarrow B \rightarrow G \quad | \quad G \rightarrow H \quad K \rightarrow L$
 $C \rightarrow D \quad | \quad I \rightarrow J \quad I \rightarrow D$
 $J \rightarrow K$

CAUSALITY IS REACHABILITY IN SPACE-TIME.



look no further

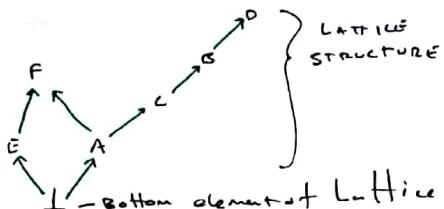
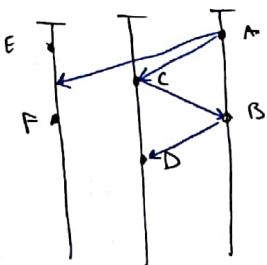
messages entangle causality.

"Happens Before" → CAUSALITY

given two events, (a, b)
a happened before b if a occurs at an earlier time than b.

- LAMPORT
- given two events, (a, b)
a happened before b if (a, b) occur in same process & a occurs first on a single process line.
- a could cause b if (a, b) occur in same process & a occurs first on a single process line.
 - if a is a send event & b is corresponding receive event
a happened before b.
 - if there exists a C, such that a happens before C & C happens before b
then a happened before b. $\exists C \text{ s.t. } a \rightarrow C \wedge C \rightarrow B \text{ i.e. say } a \rightarrow b.$

PARTIAL ORDERS



1. TRANSITIVE $(a \rightarrow b, b \rightarrow c) \Rightarrow (a \rightarrow c)$
2. irreflexive $a \not\rightarrow a$ (Asymmetric)
3. Asymmetric irreflexive $\neg a \rightarrow a$

↳ bot. element / child of lattice such that no element is less than.
↳ Diagram to catch all ordered pairs.

TOTAL ORDERING

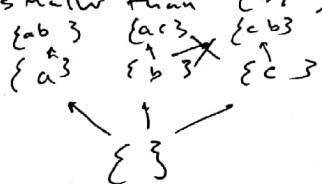
PARTIAL ORDERING w/ EXTRN RULE

- for all (a, b) $a \rightarrow b \vee b \rightarrow a$

PARTIAL ORDERING: (Subset example)

$\{A, B, C\}$ consider all subsets of powerset (A, B, C)

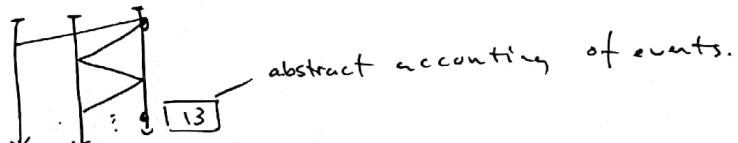
$\{B\}$ is smaller than $\{B, C\}$ because b is contained in set $\{B, C\}$



elements are comparable if a path through the lattice exists.

Logical Clock

- 1) Lamport clock
- 2) vector clock

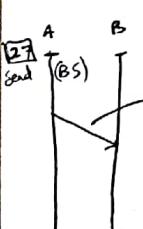


A) Mutable **state** containing clock value.

e.g. Lamport clock \rightarrow integer

vector of integers.

B) **Interpose** - get between-on communication, such that, all messages contain clock value.



message sent contains payload & metadata (clock) (current estimate # relative events)

C) SET OF RULES DESCRIBING HOW & WHEN TO MODIFY CLOCK.

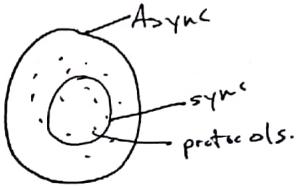
1. $\rightarrow (a \Rightarrow a)$
 2. $(a \Rightarrow b) \nRightarrow b \Rightarrow a$
 3. $a \Rightarrow b \wedge b \Rightarrow c \Rightarrow a \Rightarrow c$
 4. $\forall a, b, a \Rightarrow b \text{ or } b \Rightarrow a$
- } TOTAL ORDER
} PARTIAL ORDER
} BEST WE CAN GET IN D.S.

Models

1. Sync \rightarrow CPS CAN'T BE ARBITRARILY SLOW
2. Async \rightarrow

Proofs

1. existence / constructive / correctness
2. Impossibility.

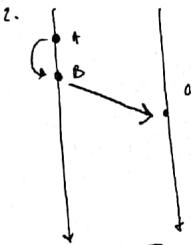


strongest correctness proof \rightarrow Async Model proves Sync model. \rightarrow Proves so in Async world, assuming worst-case.
prove impossibility. \rightarrow do so in a best case Scenario

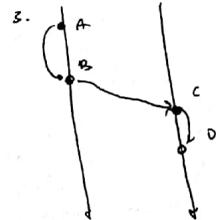
$\rightarrow 2, 2 \Rightarrow 3, 1 \Rightarrow 3, 2 \Rightarrow 4, 1 \Rightarrow 4$

Happens Before A "happened before" B = $A \Rightarrow B$

1. A, B in same process, A Happens before B.



"causality is reachability in Space-time"
Cause of 4? \rightarrow dots that can reach 4



Concurrency
if $a \Rightarrow b \wedge b \Rightarrow a$
true
 $a \parallel b$
 a is concurrent to b
"no causal / relation"
when two events are concurrent, they can float
down lines freely.
it $a \parallel b \wedge b \parallel c \neq a \parallel c$

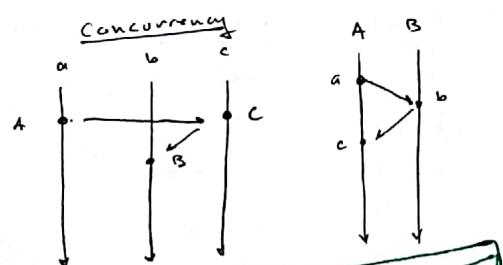
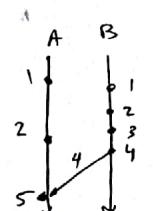
Import Clock Rules

1. Every Process maintains an integer counter initialized to 0.

2. Local Events increment process counter.

3. Processes attach counter to every message.

4. When processes receive a message, they update their counter to $\max(\text{local counter}, \text{msg counter}) + 1$



"consistent with causality"
compare: $<$ our integers
merge: $\max()$ our integers

VECTOR CLOCKS

$$a \rightarrow b \iff V(a) < V(b)$$

↳ stronger than Lamport Clock
"characterize causality"
equivalent to happens before.

- 1) $a \rightarrow b$
- 2) $b \rightarrow a$
- 3) $a \parallel b$

Compare

$$\begin{bmatrix} 1 & 2 \end{bmatrix} < \begin{bmatrix} 2 & 2 \end{bmatrix}$$

↓
later in "causal time"

$$> \begin{bmatrix} 0 & 0 & 1 \end{bmatrix}$$

$$\parallel \begin{bmatrix} 2 & 1 \end{bmatrix}$$

compare: $V(a) < V(b)$ if $\forall i: V(a_i) \leq V(b_i)$ every element in a is greater than or equal to every element in b .
Merge:

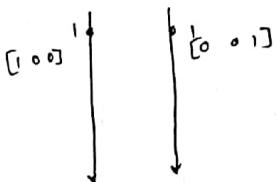
3 numbers

$$\begin{bmatrix} A & B & C \end{bmatrix}$$

↳ represents (integer value) events committed on each timeline.

$$V(a) < V(b) \text{ iff } \forall i: V(a_i) \leq V(b_i)$$

$$\exists j: V(a_j) < V(b_j)$$



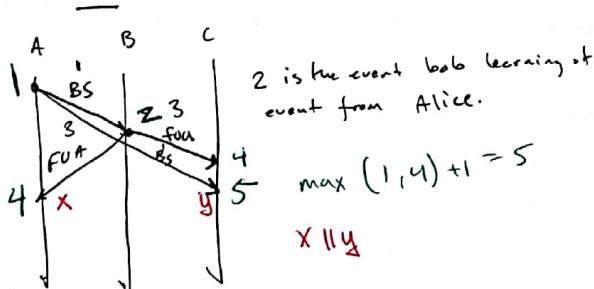
Merge

ELEMENT-WISE MAX.

VC

1. every vector is a vector of ints
2. increment "hour" integer
3. processes attach counter to every message
4. ELEMENT-WISE MAX for updating

LC



$$\max(1, 4) + 1 = 5$$

$X \parallel Y$

GREAT GUARANTEE.

CAUSAL TOTAL ORDER

FIFO TOTAL ORDER

DELIVERY PROTOCOLS

TOTAL ORDER

FIFO

CAUSAL

weak causal

strong causal

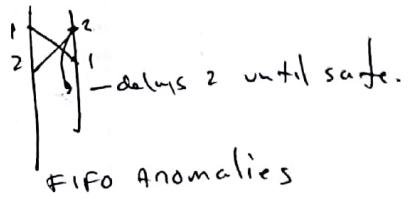
IF A SENT m_1 before m_2
B DELIVERS m_1 before m_2

CAUSAL D

IF $A \rightarrow B$

A IS DELIVERED BEFORE B.

Any FIFO violation is a violation of causal, but no opposite.
CAUSAL is STRONGER THAN FIFO.



FIFO anomalies

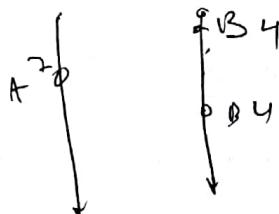
CAUSALITY -

Any pair of events are comparable.

$$C_a = 7, C_b = 4$$



a doesn't happen before b
is anything, we can conclude.



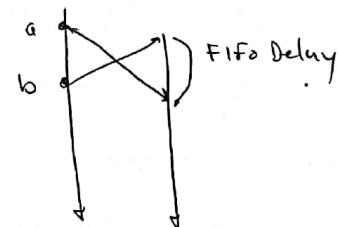
ORDERING GUARANTEES (CATLUS)

FIFO - network acts as a queue
 $0 \rightarrow 0$ one-to-one

broadcast to many

many - to one

many - to - many

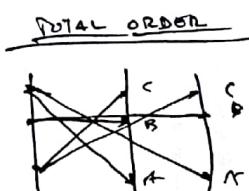
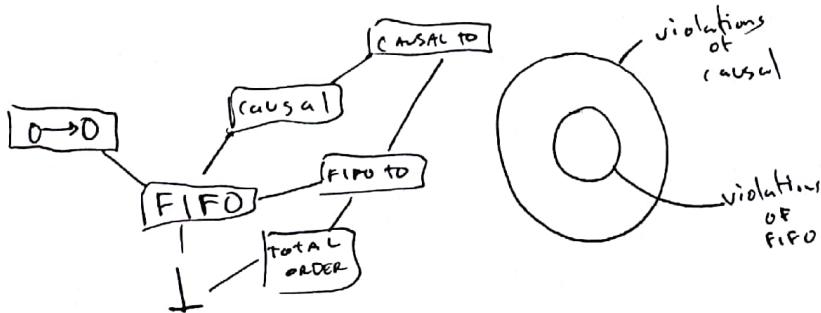
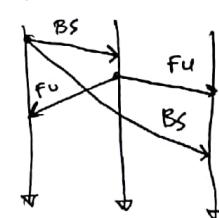


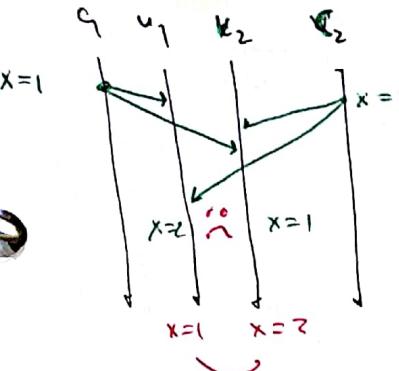
CAUSAL = FIFO (violation)

↳ violations of causal ~~and~~ but are not violations of FIFO, because causal is stronger



causal violation





So, this is a violation
of total order.

not causal PO violation, because a & b are concurrent.

if a happened before b
deliver a before b.

should be.

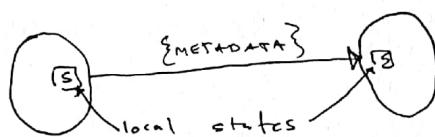
DELIVERY PROTOCOLS

• A SET OF RULES

→ COMMUNICATION → NEED MORE THAN 1 PERSON TO NEED A PROTOCOL.

→ CONTRACT BETWEEN COMMUNICATING PROCESSES.

PAIR OF COMM. ENTITIES



DELIVERY ORDER PROTOCOLS

1. LOCAL STATE.

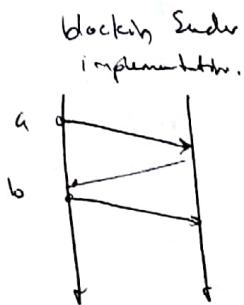
2. RULES FOR SENDING MESSAGES

3. RULES FOR RECEIVING MESSAGES

4. WHEN IT'S OK TO DELIVER A MESSAGE

→ MIGHT RECEIVE MESSAGE, BUT
DELIVERY CAN BE DELAYED.

IN FACT, THIS IS ALL D.O.P,
D.O.



don't send b until
ack, a has been
received.

causal
Send by: everything happens in thread.

autonomous vehicle
comm. via Radio.

non-blocking algorithm

- can send messages w/o waiting
- RECEIVER sorts out FIFO.
- senders have counter
to include
- receiver has counter for every sender
delivers message if it is expected
- eventually all messages are FIFO delivered.

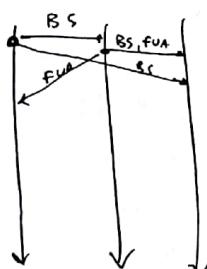
State:

Sender: Every message I have ever sent
receives B

~~rules~~:

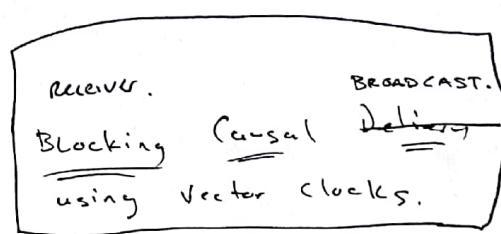
rules: When sending a message, attach all messages b , such that $a \rightarrow b$.

(2)

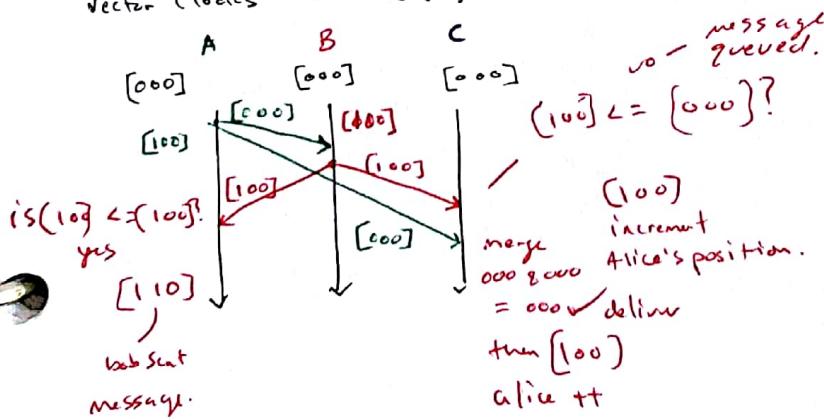


uniquely identify inputs & remember.

BROADCAST REGIME



vector clocks characterise causality.



STATE:

sender: V.C.
receiver: V.C., queue

rules:

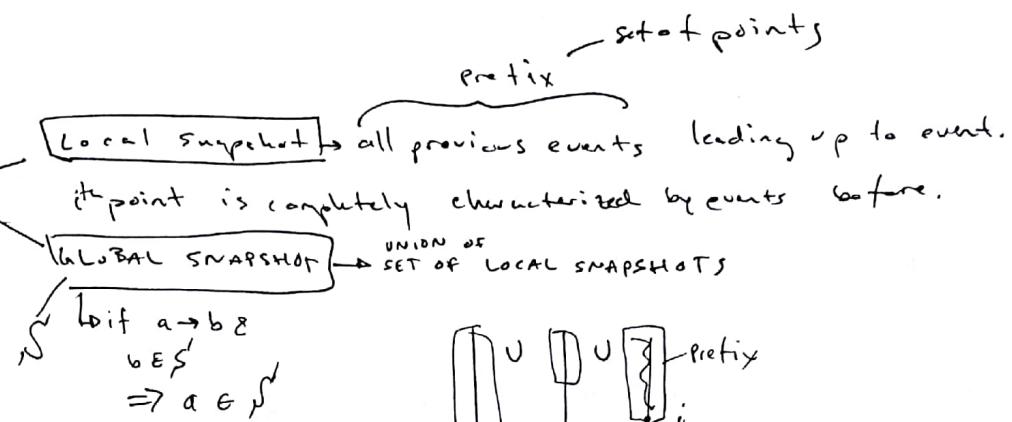
1. after message send, increment my position in V.C.
2. on receipt, merge V.C. w/ local V.C., then increment senders position.
3. on send, include local V.C.
4. Deliver message with V.C. V when $V \leq V_{local}$

n positions in V.C.

Omission Model

- DROPPED MESSAGES
- ACT AS CRASH

CONSISTENT GLOBAL SNAPSHOTs



not reasonable

ALGORITHM FOR OBTAINING GLOBAL SNAPSHOT

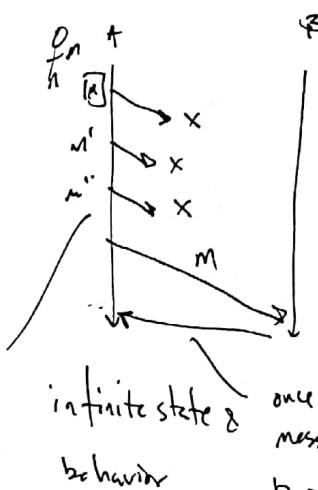
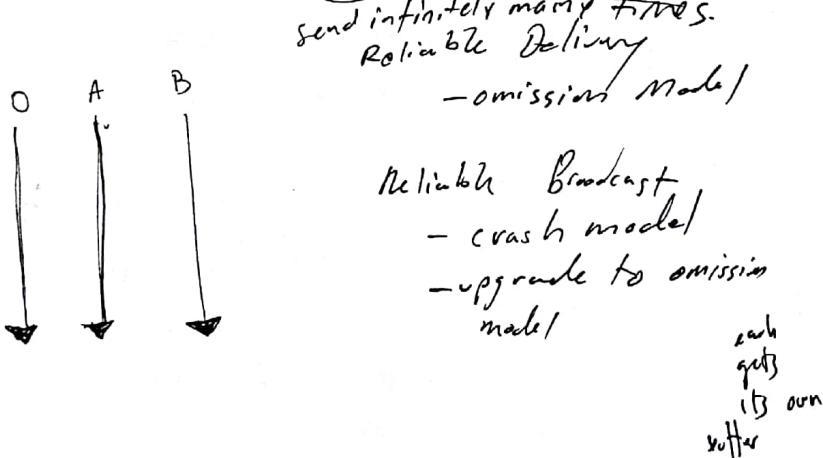
Chandy / Lamport Snapshots

can be an external witness

1. An observer process sends a token to every other process
2. When a process receives a token, it records its local state, & sends to observer, then sends a token to all processes, ^{for the first time} & begins recording in-bound messages.

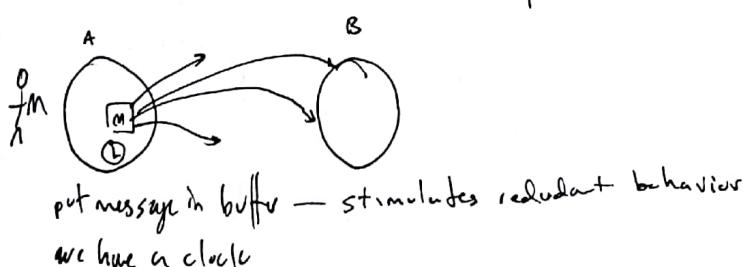
ASSUMPTIONS

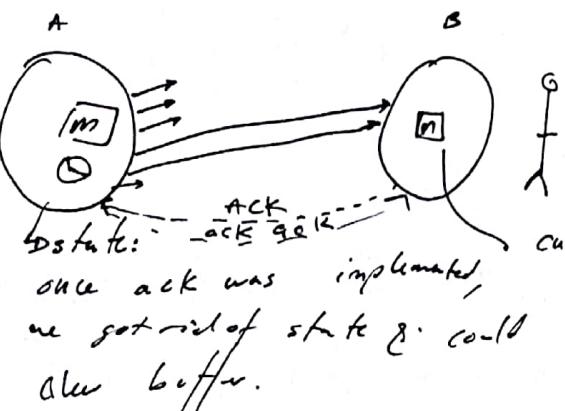
1. NO FAULTS
2. FIFO DELIVERY
3. ALL-TO-ALL COMMUNICATION



NOTE: everyone sends and receives a token from everyone.

Tokens always push relevant messages in snapshot.
 $B \rightarrow A$
 saying M was received, A's local state changes.





now, ACK ACK clears up B's state.

Exactly once delivery algorithm process

can have buffer on receiving end.

eventually, if not ever, if not all messages are dropped, RD will work.

2^{64} bit number

for sequencing
is how TCP

handles two general

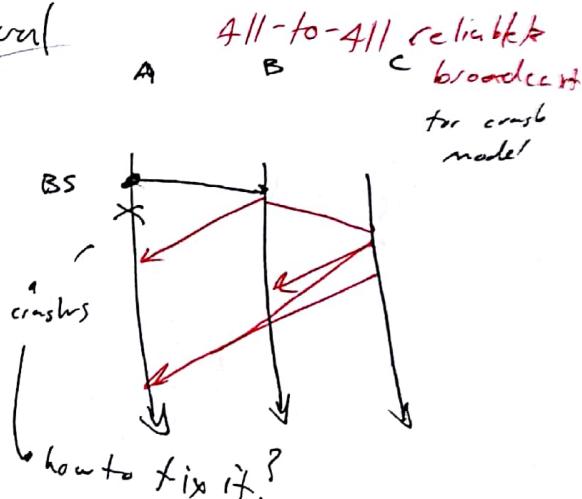
Reliable broadcast
- crash mode /

if a correct process delivers a message, then all processes deliver it.

if A then B $A \Rightarrow B$

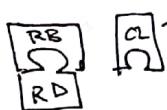
$\neg A$

$A \wedge B$



↳ when a process receives a message, it relays the message to every process it can touch.

(3)



Independent to each other wrt FAULT

FAULT TOLERANT Protocols

FAULT TOLERANCE IS REDUNDANCY

Alt

SOME FAULTS WON'T AFFECT ALL RESULTS

TIME: RETRY

SPACE: REPPLICATION

MATH: Error-Correcting Codes.

Representation: DATA PAGES / write-through / SHADOW PAGES.

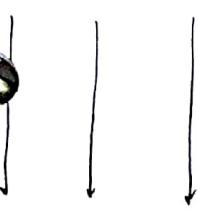
Program: N-version Programming → LIFE CRITICAL SYSTEM.

Two copies of state asynchronously updated can raise caution about synchronization & which data values are correct.

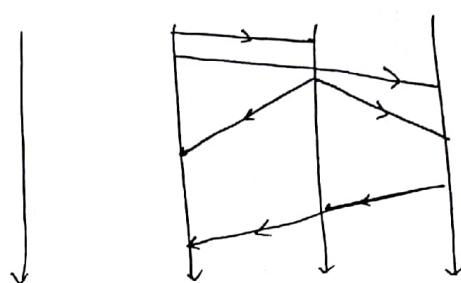
- Assume INDEPENDENCE
- REDUNDANT BEHAVIOR, COMPUTATION, STATE, etc.

SECTION

0 A B

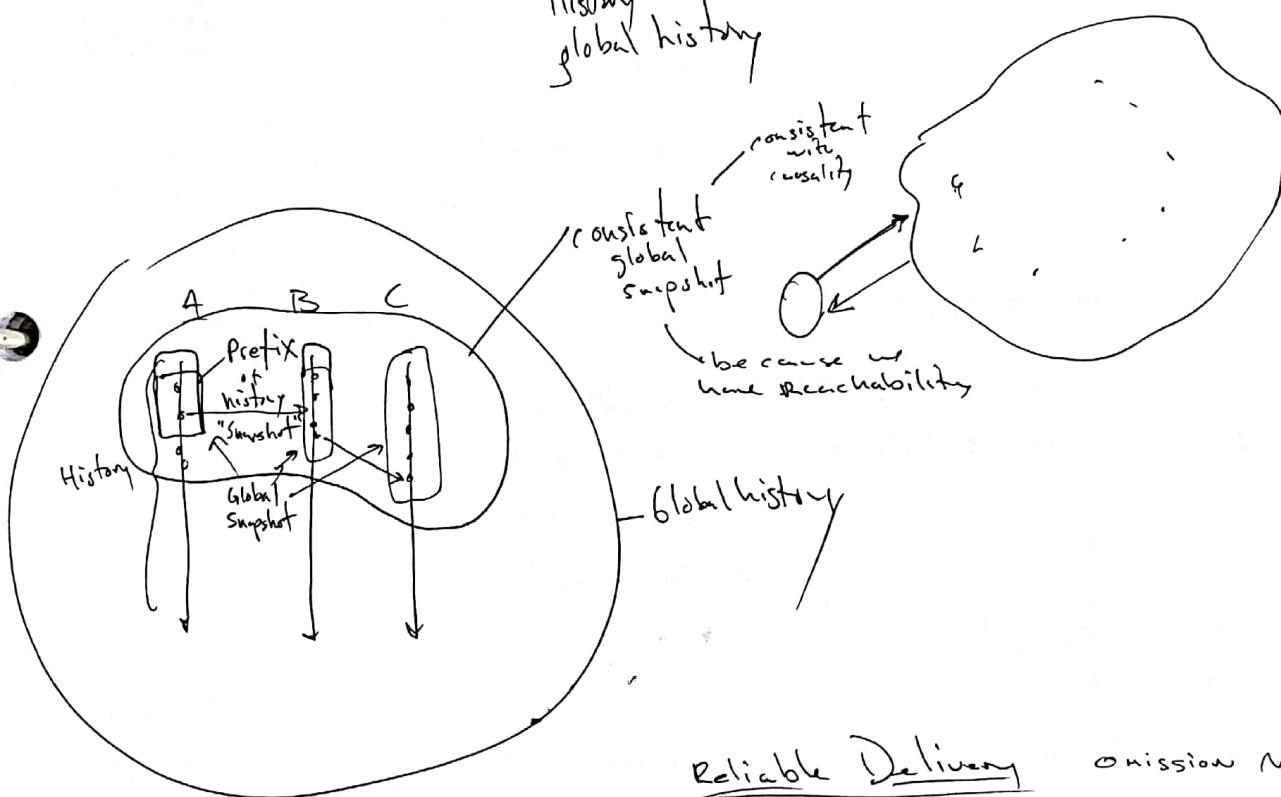


0 A B C



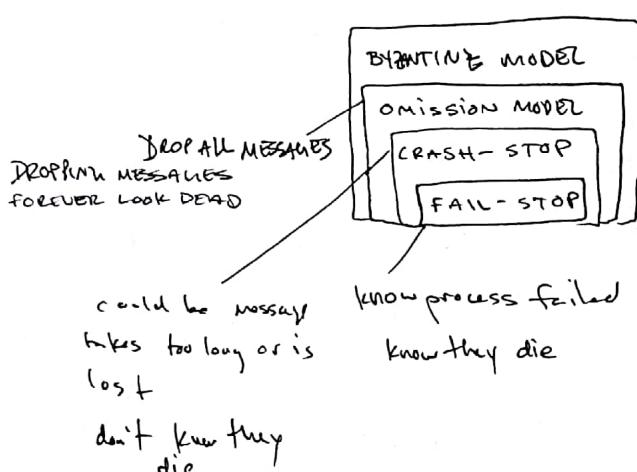
each has
a single queue
for each process
other than itself.

VC characterize causality

History
global historyReliable Delivery

omission Model

FAIL-STOP MODEL



DELIVERY : eat cake

SEND : PASS cake

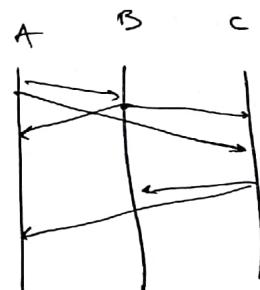
RELIABLE BROADCAST

if one correct process delivers, all do. will send message.

Send first, then open

Short task? later.

use RELIABILITY DELIVERY

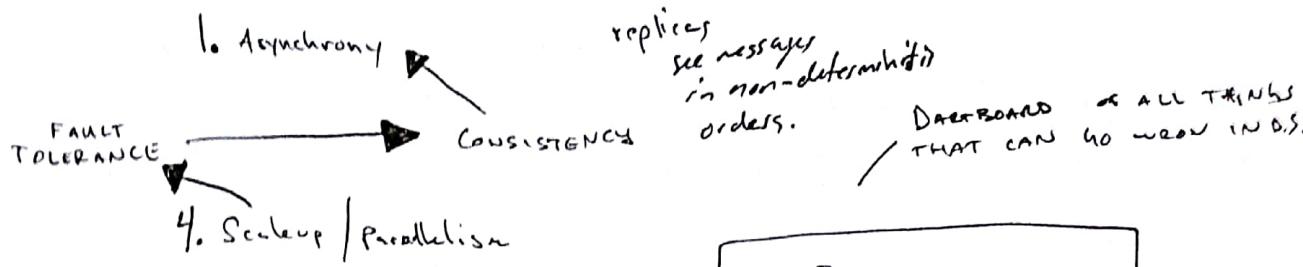


DYNAMO

FAULT-TOLERANCE is REDUNDANCY.

Time Partial-failure in space

Asynchrony - message doesn't arrive "in time" - keep sending



$$P_{fail} = \frac{1}{1000}$$

if we scaling to 1,000 computers
we start to approach the apex
of the parabola.

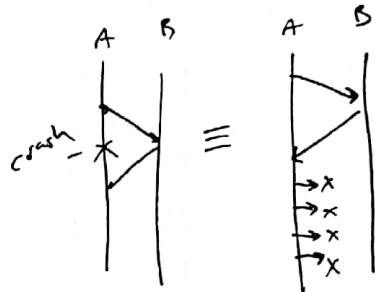
$$1000p = 1 = f(n)$$

well, if we have 3 cpus
acting in a D.S. as 1/cpu.

we get $\frac{3000}{4 \text{ lines}} \cdot \frac{p^3}{\text{scale}} =$
exponentially better
probability of failure.

MESSAGES
DON'T GET
DROPPED,
SERVICES
FAIL.

- CAN BE MODELED BY SUBTRACTION EVENTS.
- REDUNDANCY IN TIME OVERCOMES OMISSION.

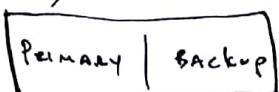


Replicate data & computation

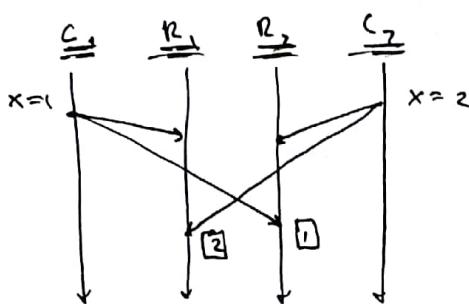
we want to go fast (scalability),
but we also want to go (fault tolerance)
so we make replicas (consistency), but we
need to make sure they're in sync, due to
unordered network (async).

- if it doesn't pass
the fluffy higher model,
it won't pass (narrower ones)

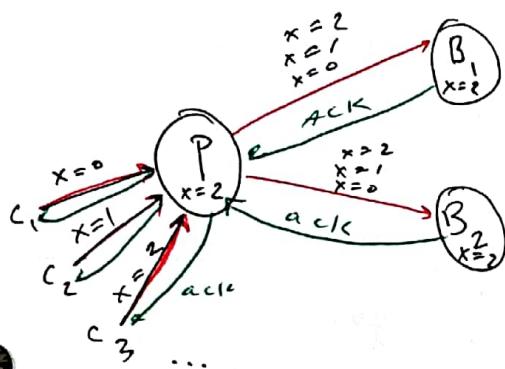
simplest.



- CAUSAL DELIVERY USING VECTOR CLOCKS.



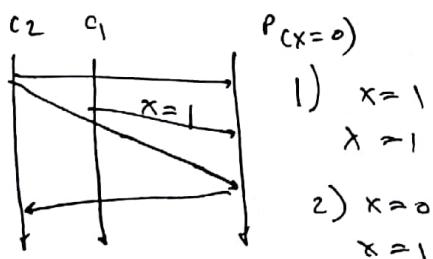
Primary - LEADER (Computer)
Backup - followers



Strong Consistency

There is no way to tell the system is distributed.

→ REPLICATED SYSTEM is S. CON.
if there can be no witness
of the fact it's replicated.



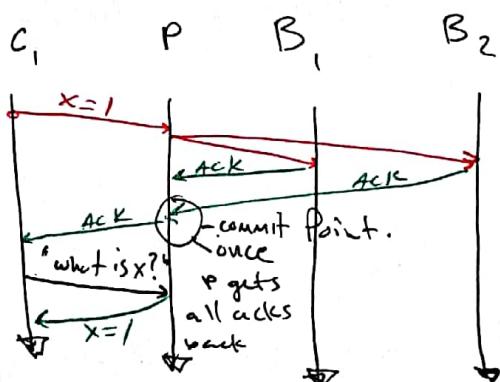
- Don't permit clients to query backup
- don't ack client until all backups ack.

$C_1 \neq R_2 \rightarrow$ inconsistency.

PROBLEM: CAUSAL ORDER NOT TOTAL.

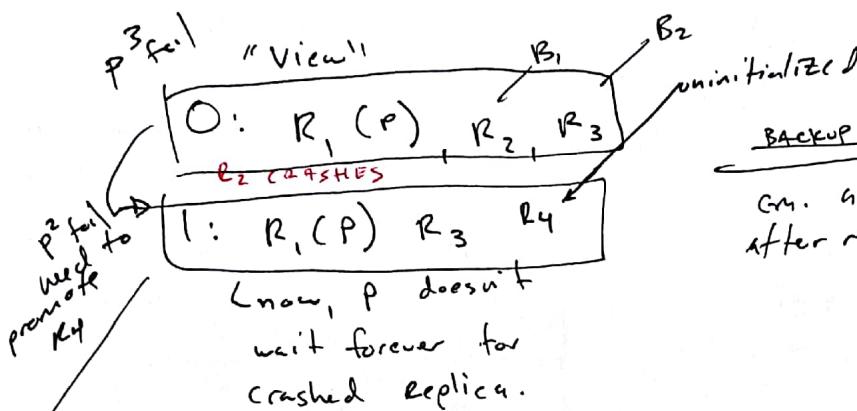
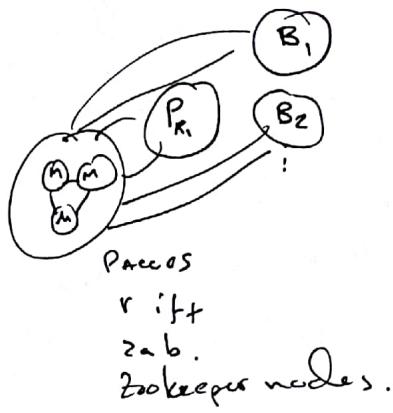
ESTABLISH TOTAL ORDER

AFTER (P) receives ACK from ALL!
backups, it will ack ~~the~~ client.



pull primary to get value of x.

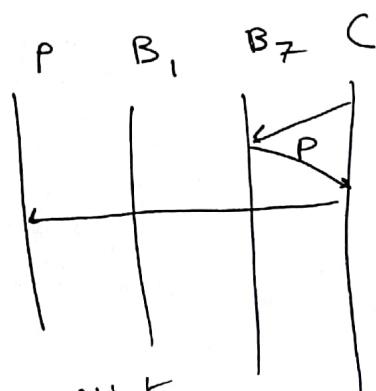
(3)

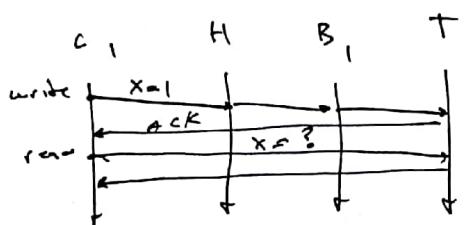
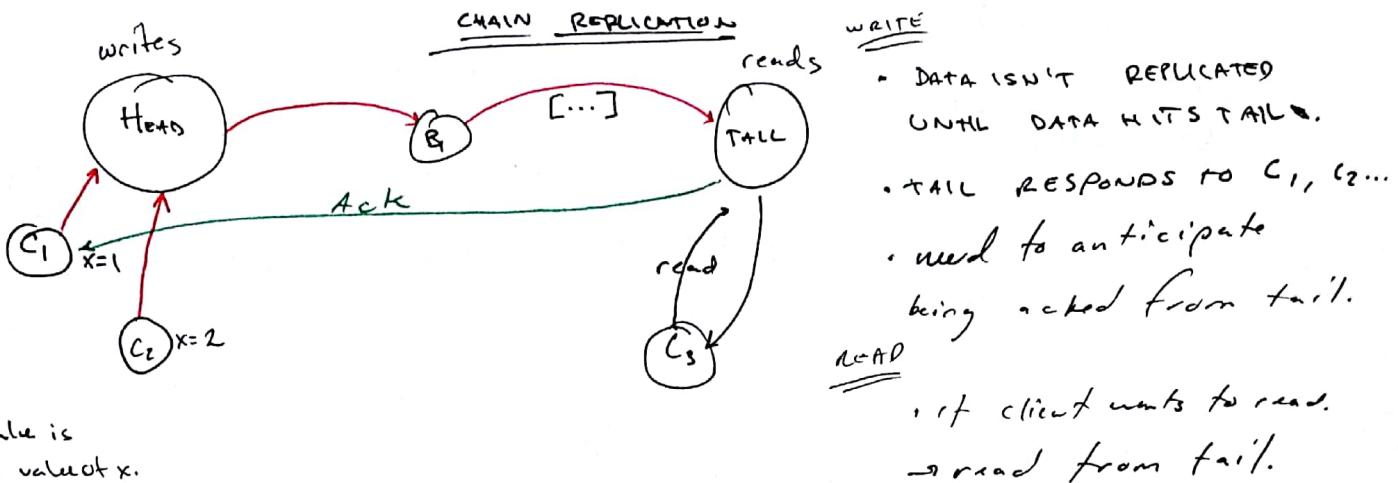
FAULT-TOLERANCEcluster of computers \rightarrow Configuration Managercluster of (M. circle), talking using
consensus algorithm.

Set of views atomically propagated to all nodes/replicas.
What happens when we crash?

want to keep $P_{fail} = P^3$.
do this by promoting already existing replica node waiting
for a time like this.
 \rightarrow note: it's uninitialized, needs to be dealt with.
 $\hookrightarrow R_4$

$2 : R_3(P) R_4 R_5$ \rightarrow view change needs to be visible to clients as well





LATENCY → client perceived latency
 Primary / Backup → Max Roundtrip time to the (slowest) backup.

Chain Replication - sum of latency of RTT's

more bytes through network.

among nodes in chain.

→ better network utilization

→ less work, proxy doesn't handle dissemination. or does.

→ Pipeline writes

→ shed the load of read to tail

RELIABILITY - likelihood of a component failing at a given time.



AVAILABILITY - likelihood a system is functioning correctly at any time.

$$A = \frac{\text{uptime} - \text{how long ran before crash}}{(\text{uptime} + \text{downtime})}$$

.9999

4 nines of availability.

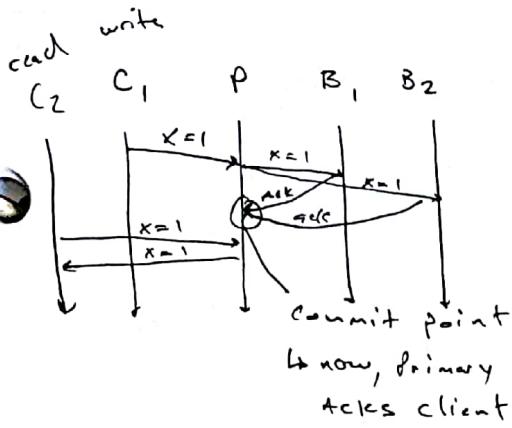
MTBF - MEAN TIME BETWEEN FAULTS

$$A = \frac{MTBF}{MTBF + MTTR}$$

MTTR - MEAN TIME TO RECOVER

To cut this, increases availability.

$$\left(\frac{1}{1 + \frac{MTTR}{MTBF}} \right)$$



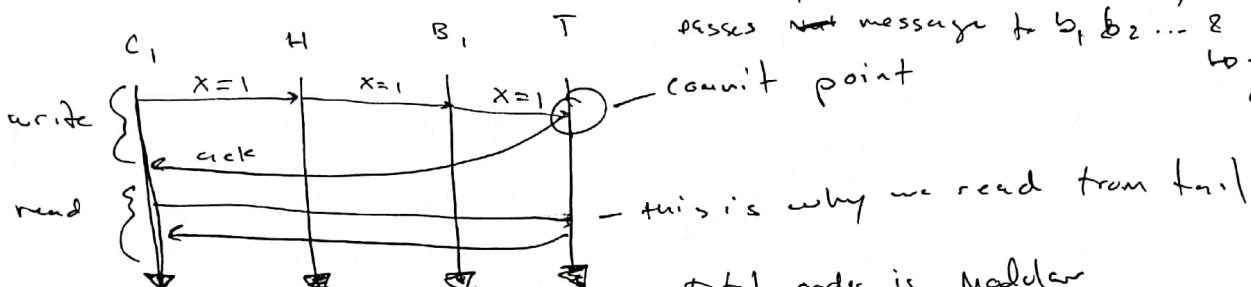
latency = max latency edge.

replicated systems

- Primary Backup
 - ↳ decides order based on arrival time
- Total order over operations

latency = sum of edges.

CHAIN REPLICATION



Total order is Modular

Head chooses order based off delivery order of writes

tail

↳ chooses order of reads.



DELIVERY Protocols - Correctness Properties

Safety. & Liveness

any correctness property is one or both.

Safety: "A bad thing never happens"

always have finite counterexamples.

Liveness: "a good thing eventually happens"
Infinite execution.

SAFETY
if you can draw it, safety property

A B S

CROSS DELIVERY - bob smells



FIFO

LIVENESS

RELIABLE DELIVERY

R.D. eventually, a message will get delivered

(2)

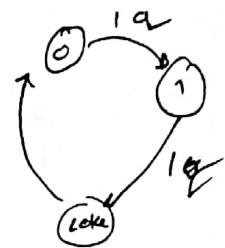
Dynamo
paperSMR - STATE MACHINE REPLICATION

1. STATE MACHINE
↳ processing is deterministic

2. Assume Ordering Oracle

$$f: (S \times I) \rightarrow S'$$

S = state
I = input



Client asks O for ordering.

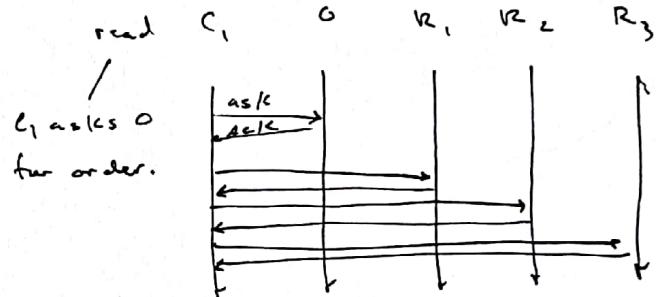
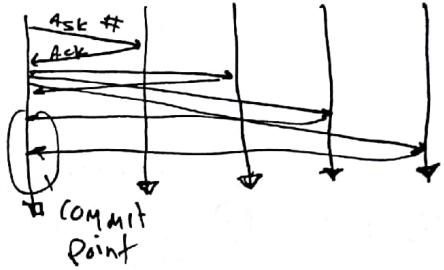
Client plays role as primary, SMR ASSUMES SMART CLIENT

client disseminates data.

O decides order.
C_i sees to it all R get data

-symmetrical

C _i	O	R ₁	R ₂	R ₃
----------------	---	----------------	----------------	----------------



implement O via consensus.

token from O to C_i gets attached to messages for everyone to compare.

1) clients send operations to a Set of replicas.

Note: Set is not necessarily all replicas.

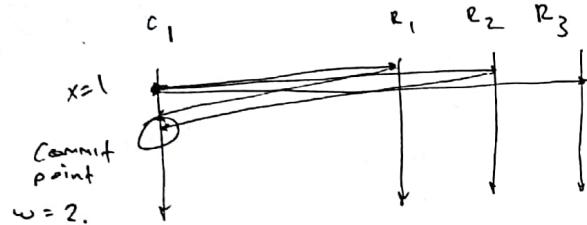
2) Three key parameters

N - number of replicas in system

w - number of acknowledgments required for a write

r - number of ACKs " " ^{read}

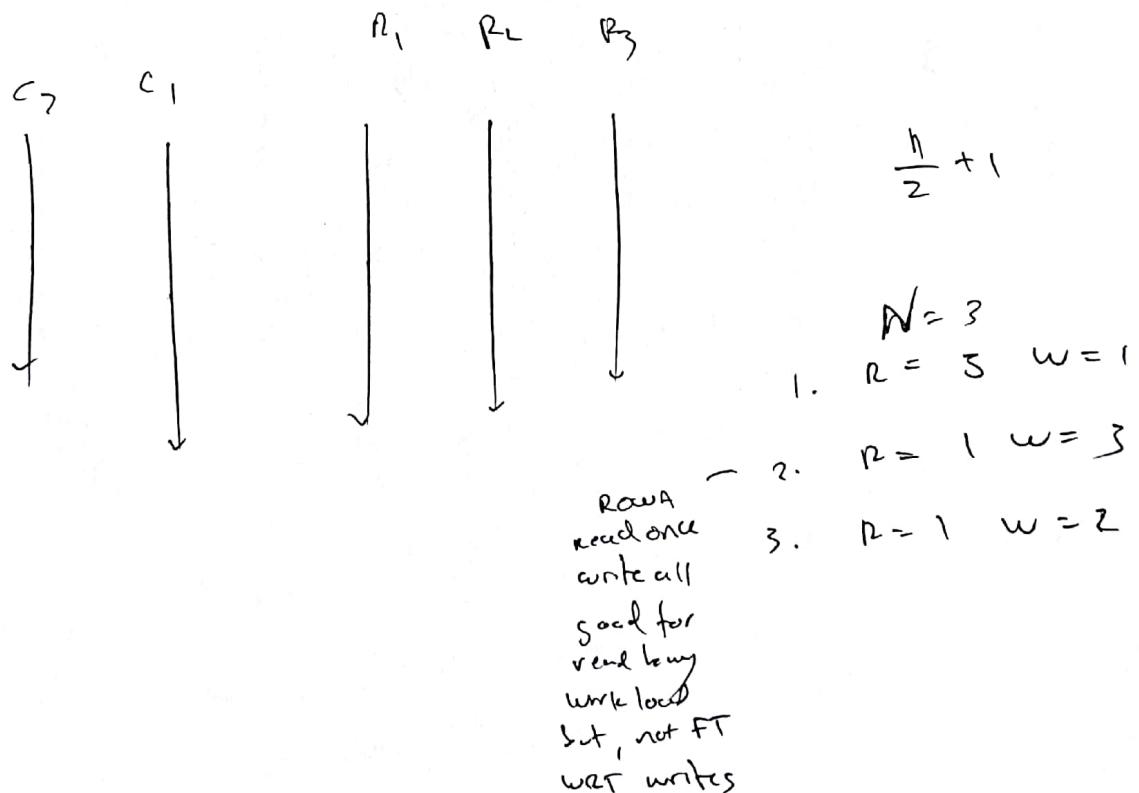
Ex: $\langle N, w, r \rangle = \langle 3, 2, 2 \rangle$ - fastest replica determines latency



FAULT TOLERANT

Strongly Consistent

$$w = \left(\frac{R}{2} + 1 \right)$$



(4)

Quorum RESTRAINTS

1. Every read quorum should intersect with every write quorum

$$R + w > N$$

2. Every 2 write quorums must intersect

$$w > \frac{N}{2}$$

equates to No concurrent writes.

CAP THEOREM

□ Consistency - single copy consistency

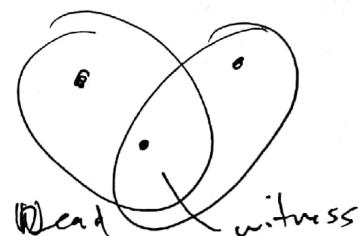
R=1
w=2

□ Available - as long as a node exists, it should provide a meaningful response

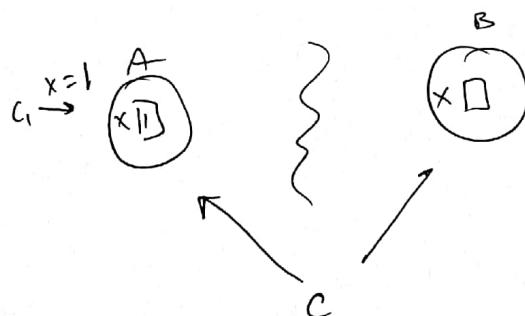
□ Partition tolerant

pick 2 out of 3

every D.S. is CP or AP. No D.S. is C.A.P.



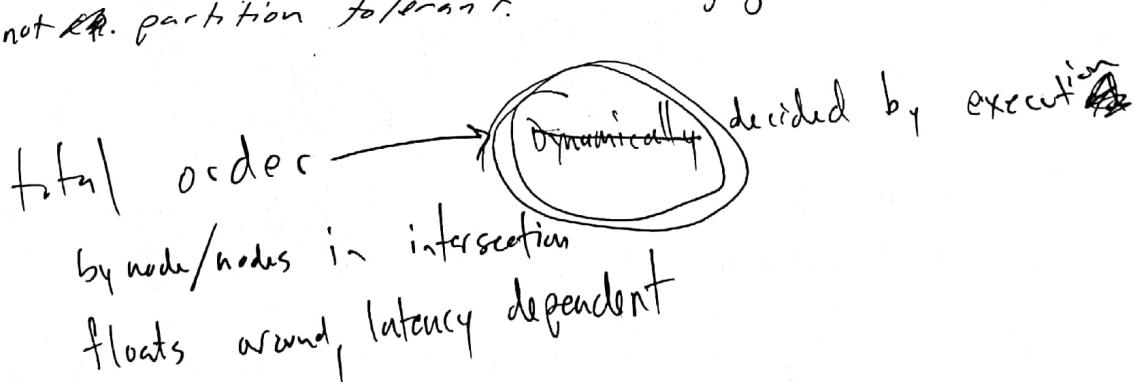
if C queries B and B responds I don't know. (CP.)
to order at writes.
in leases.



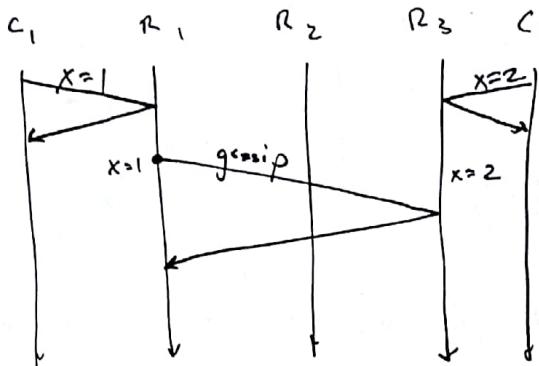
Satellite - AP

any situation where we have majority quorums → not CP. partition tolerant.

google

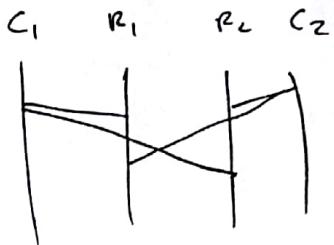


Inconsistent gossiper replication



$$\begin{array}{l} w=1 \\ r=1 \end{array}$$

FAST!



need Vector Clocks to handle gossip protocol.

can't rely on clients to follow protocol.

it disjoint,

merge

it not,

~~topic breaker~~

gossip
- pairwise
↳ Anti-entropy

~~PERIODIC~~
↳ Done by V.C.
concurrent

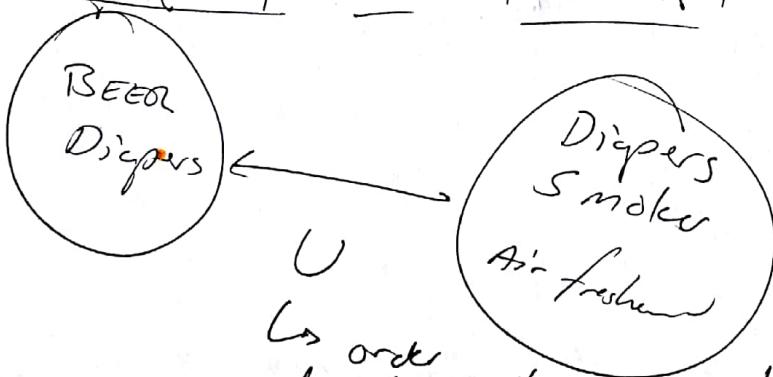
(data goes both ways
~~not~~ one process has a lot of data over the other. &

(data goes one way

Total order anomaly
FIFO CORRECT

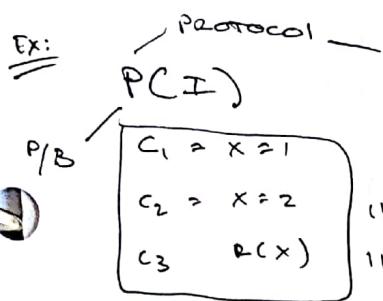
$f(x) = f(f(x) \circ f(f(x)))$
Commutativity
↳ can retry with impunity.

Shopping Cart Analogy

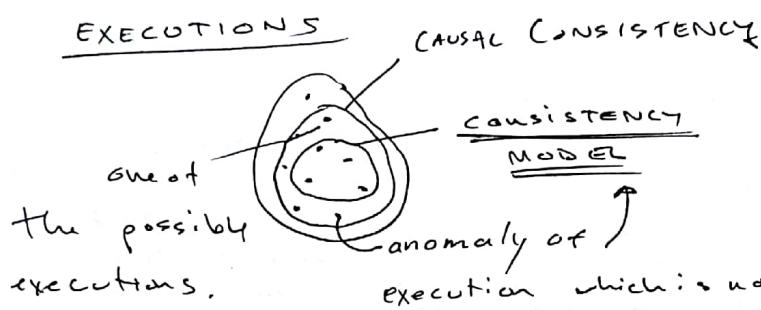
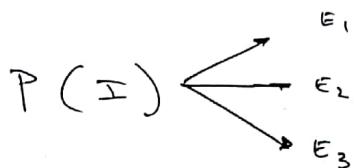


↳ order doesn't matter due to commutative property of sets.

Set deletion doesn't have the ↗



a given protocol gives rise to
Set of EXECUTIONS



LC's ARE CONSISTENT WITH CAUSALITY

$$a \rightarrow b \quad LC(a) \subset LC(b)$$

$$b \not\rightarrow a$$

VCS characterize causality "Potential"

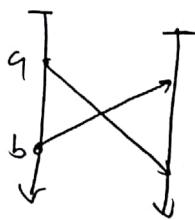
$$a \rightarrow b \iff VCS(a) \subset VCS(b) \text{ why?}$$

or

$$b \rightarrow a$$

$$a \parallel b$$

FIFO



$a \rightarrow b$

but b is delivered

FIFO total

FIFO

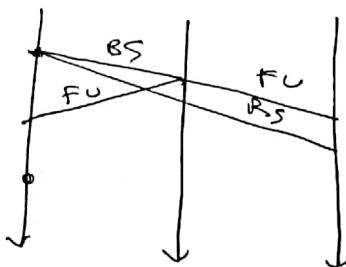
CAUSAL

CAUSAL TOTAL

Anomaly

CAUSAL

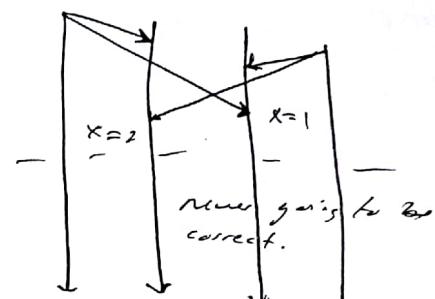
Bob Smells



CAUSAL violation
not FIFO



TOTAL



split brain
T.O. violation

FIFO allows

more executions

any FIFO anomaly
is ~~a cause~~ a cause

anomaly because
FIFO has all

executions causal
has, plus more

①

- ✓ POWERFUL PRIMITIVE
 - consensus
 - Atomic Commit
 - LEADER Election
 - Totally-ordered broadcast
 - TWO-PHASE COMMIT
- 17/7/17

SAME problem

Agreement Protocols

1. TWO GENERALS PROBLEM
2. LEADER ELECTION → FAILED PRIMARY, NEED TO BOOST A NEW NODE
3. FLIGHT CONTROL
4. DECIDING A TOTAL ORDER
5. CHOOSING A PRESIDENT
6. FAILURE DETECTION
7. GROUP MEMBERSHIP → CONFIGURATION MANAGER
8. ATOMIC CHANGES ACROSS MULTIPLE NODES



what if T is
money? Don't
want drop, or
duplicated messages.



CASE 1 - SUCCESS



CASE 2 - fail



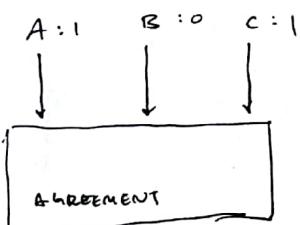
CASE 3 - fail

AGREEMENT

VALIDITY → "If a node decides x, then x was proposed"
Safety

AGREEMENT → "If a node decides x, then all nodes decide x" *most important*
safety

TERMINATION → "If a node proposes x, eventually some value is decided"
Liveness



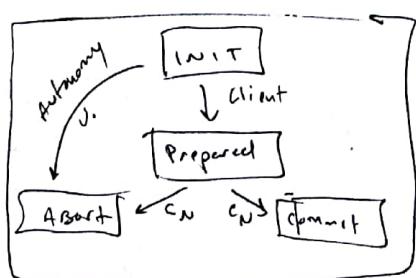
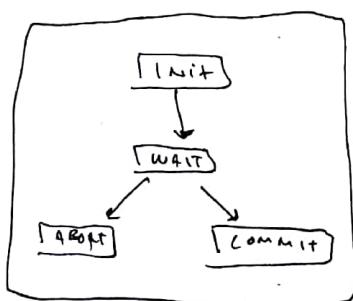
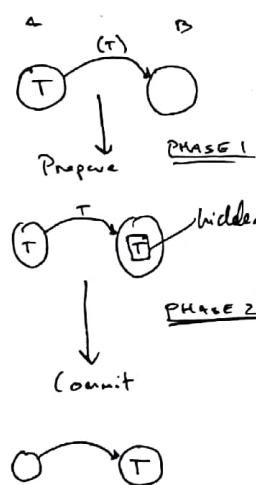
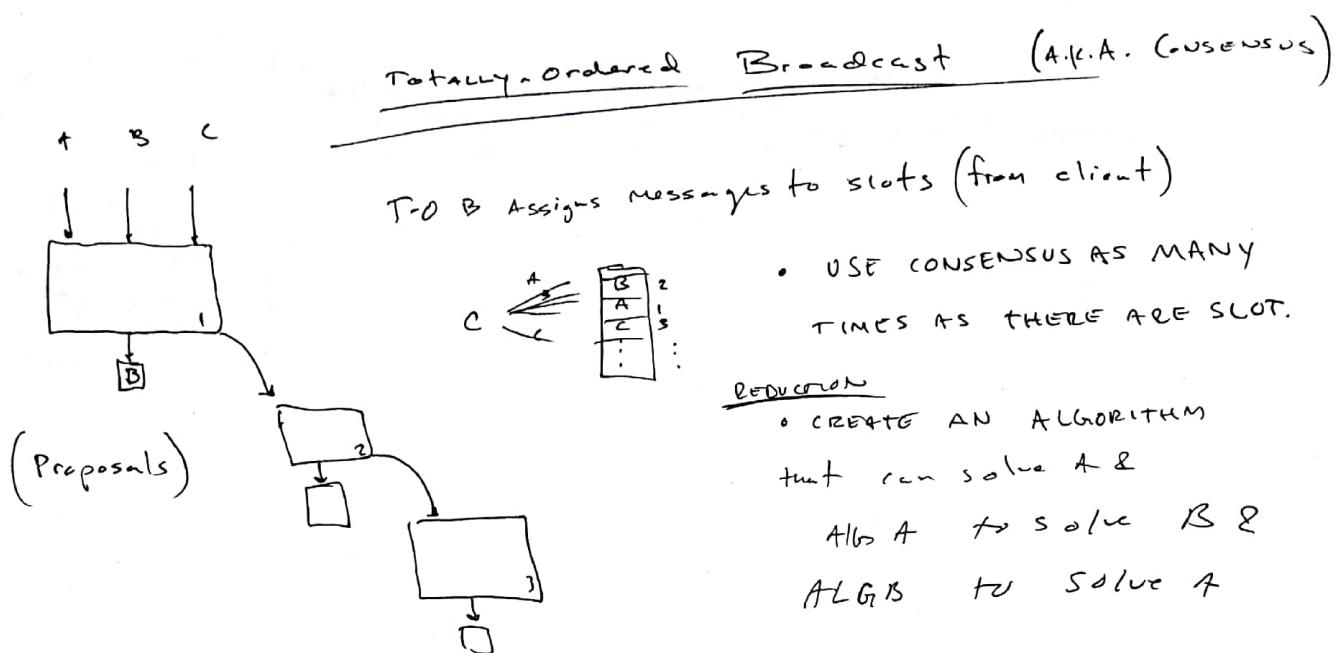
— Agreement has to
be one of the
proposed messages.

- Consider 3 instances of the above, x, y, z. All 3 instances need to agree on the output.

NOTE: Validity
is non-triviality

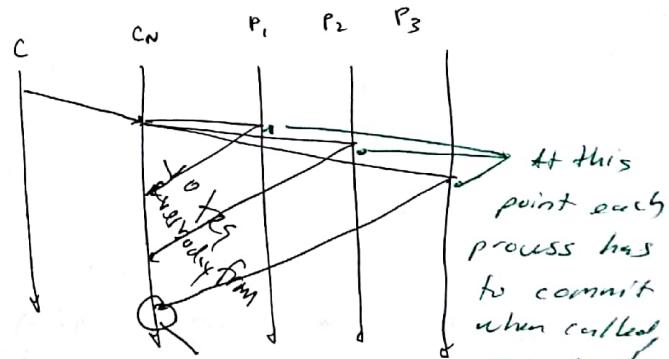
• Similar to
CAP theorem,
we typically
choose two.
Validity &
Agreement.

(2)

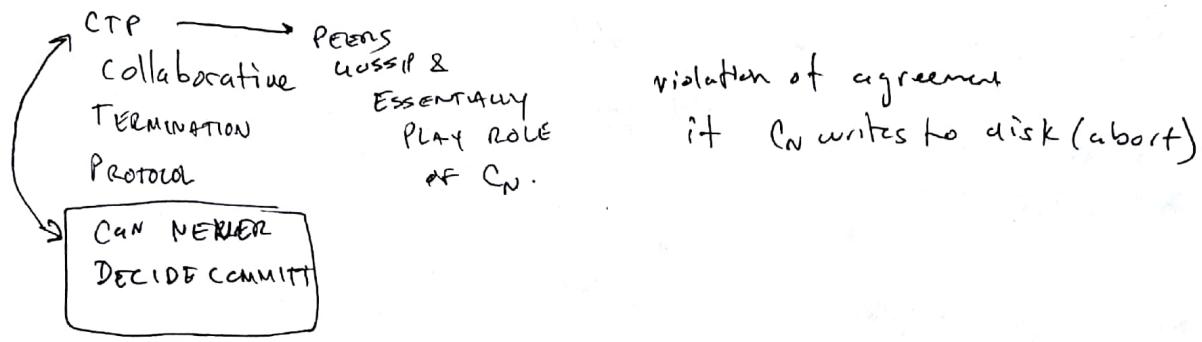


think 2 people getting wed & an authority saying I pronounce you...

coordinator node = CN - has clock.



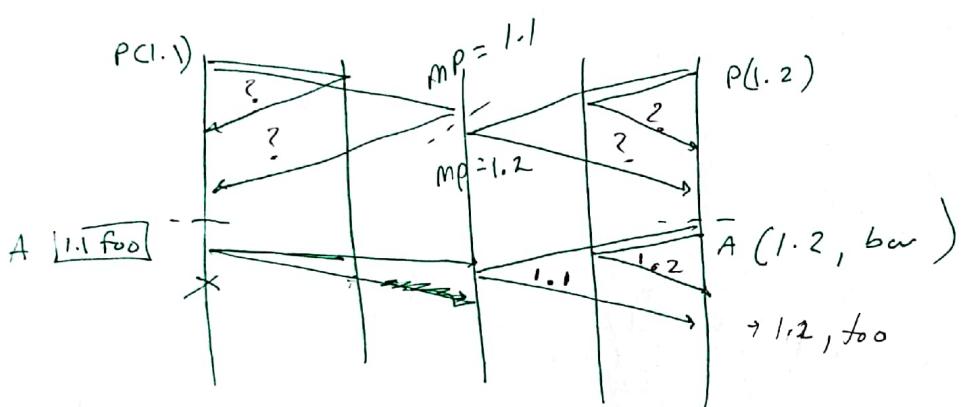
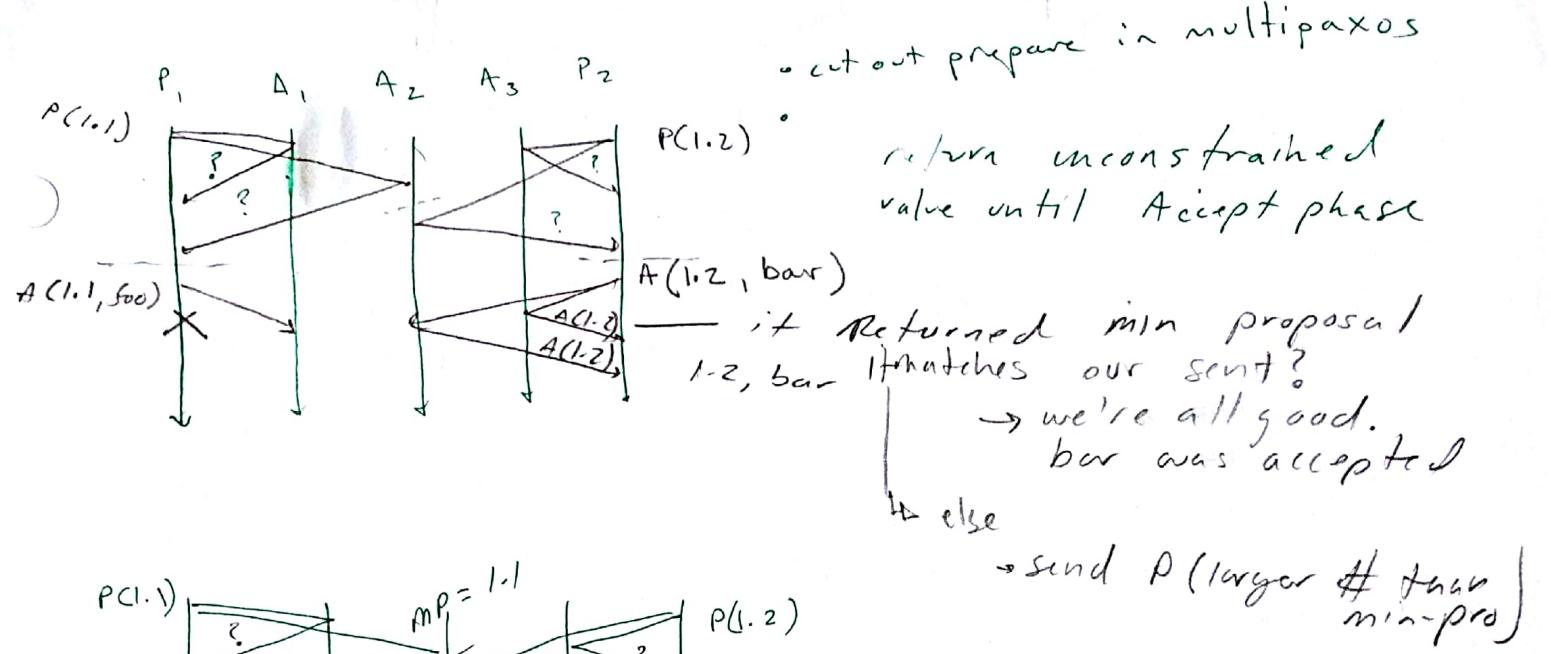
- PARTICIPANT CAN GO FROM PROPOSE TO ABORT.
- COMMIT GETS SENT ON UNANIMOUS VOTE (OR NOT) CN CHOOSES IT IS SOLELY ABOUT THE PARTICIPANTS BEING READY TO COMMIT AFTER RESPONDING.
- SEEMS LIKE PRIMARY BACK UP.
• IF A NODE CRASHES, WE COULD ABORT.



Pick a node at random &
 Start to gossip, centrally system
 will be good.

~~If in ~~the~~~~

~~Coatt → check if~~
 rather than set value to none,
 let's keep a bool dict
 tracking if we can talk to
 to nodes or not...



write-write conflict
2 xacts overwrite
Same key

write-read; exact write
variable that's
inter exact reads.

des P

non-serializable Schedules

- histories
 - invalid
 - shuffling
 - cycles (wr, wr)
 - stale read
 - divergence

- histories
 - magical
 - Shuffling
 - cycles (wr, wr)
 - stale real
 - divergence

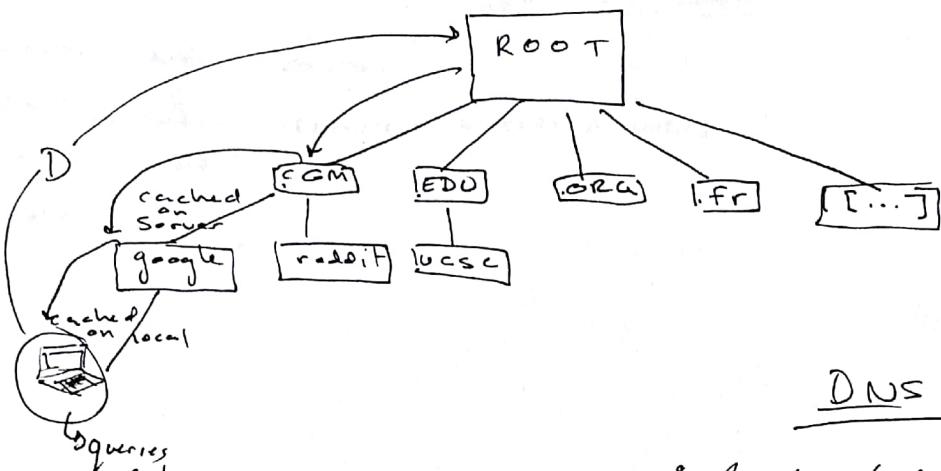
- one stream execution
- everything in order
- $x = 100$

Schedule list of transaction

list of options transaction (operations)

$w)$ if (r/w)

$$of (r/\omega)$$



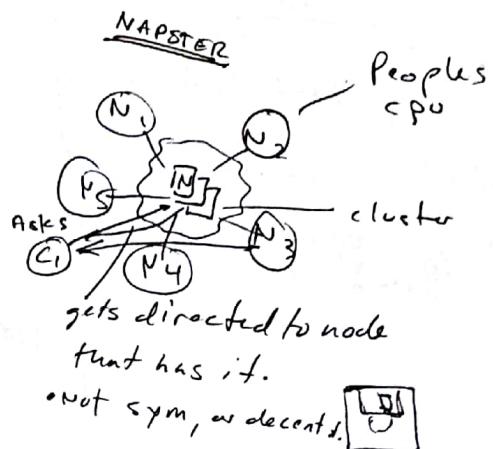
queries
Root 8
Answer/response
trickles down

DNS ENTRIES:

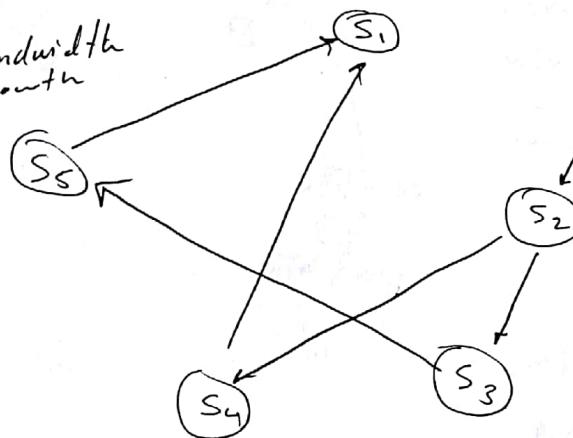
- A name place ttl load balancing service discovery
 - ↳ IPs $\{128.76.2.112, 178.76.2.113, 192.2.2.10\}$
- TTL (aggressive Caching)
 - ↳ short ttl - cache clears frequently.
 - ↳ if updating often small ttl is better, less stale data.
- each domain server decides ttl for data endorsement

PEER TO PEER

- USES, NOT SERVERS
- NO DISTINCTION BETWEEN CLIENTS & NODES
- DECENTRALIZED }
- SYMMETRICAL }

SCALEABILITY:

- minimize cost
 - Storage & bandwidth
 - want log. growth

GNUTELLA

- Decentralized
 - Symmetrical
 - flooding retrieval
 - ↳ do you have this?
- who can you see?
eventually, the
chaffing will
result in everyone
knowing.

GOALS: (weakly consistent)

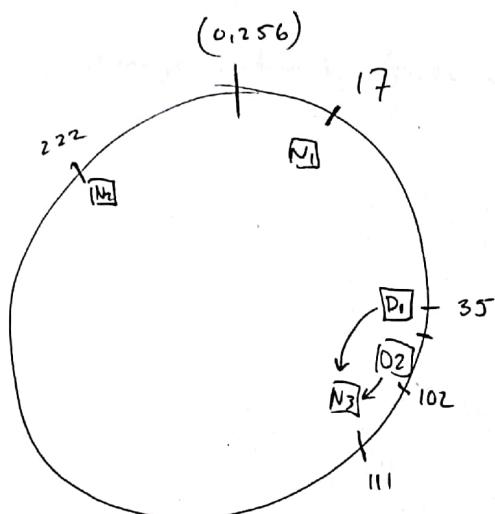
- massive Scale (100k nodes)
- 10⁶ nodes

- RESILIENT TO CHURN

CHORD - DHT

map names to nodes (computational Resources)

- SCALEABLE (Logarithmic scale storage & network costs)
- Efficient (name \rightarrow node in log time)
"lookups"



$h \Rightarrow$ hash function

$h(\text{data}) = \text{hash}$

$h(\text{node}) = \text{hash}$

Imagine we have 3 nodes

$$h(N_1) = 17$$

$$h(N_2) = 222$$

$$h(N_3) = 111$$

$$h(D_1) = 35$$

NOTE: if $h(N_4) = 45$, D_1 would range to N_4 .

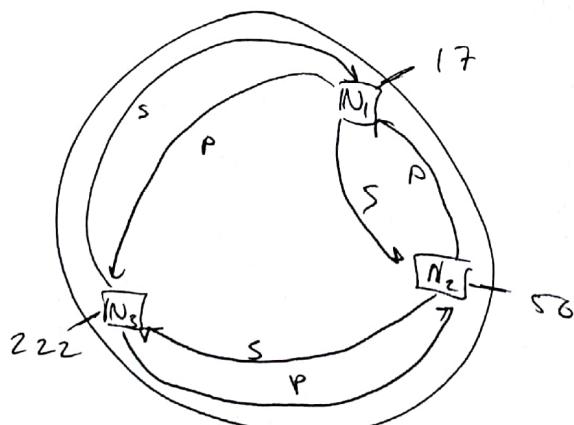
DATA ITEM BECOMES TO SUCCESSOR HASH VALUE.

$$h(D_1) = 35 \rightarrow \text{hashes/links to } N_3 = 111$$

N : NODES

K : Keys

w.H.P. if a node disappears, need to re-map $\frac{K}{N}$



STORAGE: $O(1) \rightarrow$ doubly linked list

NETWORK: $O(n)$ — have to go all around circle to find place

EFFICIENCY: Dog slow

NODE STATE	
------------	--

maintain ring { Successor (next r)
Predecessor

FINGER TABLE			
(i)	Power	Reach	ADDRESS
0	$n+2^i$		10.0.0.1
1			
2			
3			
4			
⋮			

FINGER TABLE - lives as $\log_2(n)$ for lookups.
Points to nodes powers at 2^i distance from me.

→ has $\log_2(n)$ storage pointers

So... storage, network, efficiency are all $\log_2(n)$.

11/10/17

correctness properties
for agreement

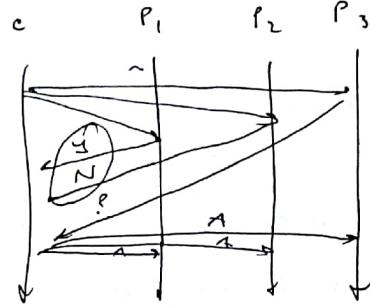
FRAGILE

TWO-PHASE COMMIT

- one-node has power
- decides → decider/coordinate
- ABORT
- C/A
- 2PC to ensure all nodes are onboard to C/A
- after votes are in,

P_1, P_2, \dots, P_n have no choice but to commit. (or abort, ideally after the votes come in yes, the $P_i - P_n$ do as C_n says)

- have to get Yes vote from every node.
- but ultimately, the Coordinator makes the C/A decision



satisfy {
1. Validity
2. Agreement
3. Termination}

CONSENSUS

1. make progress if only a bare-majority of nodes are correct.
2. To tolerate f -failures/faults, we need $\boxed{2f+1}$ nodes.
in collection of nodes, satisfy 3 agreements (well, 2 & mostly 3) and the two rules above.

(2)

Paxos

All of the computers play all of the roles

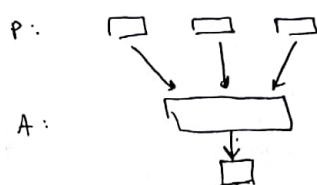
When is consensus hard?
when two opposing ideas are trying to be decided at the same time.

1. Proposer \rightarrow node \rightarrow clients may communicate \rightarrow looks similar to coordinator in 2PC
2. Acceptor \rightarrow Passive (relatively) \rightarrow listen to accept values, make decisions, send votes back regarding decision to accept values.
3. Listener \rightarrow remember which decisions were chosen.

ex: Simplest Algorithm

one acceptor

could be many proposers



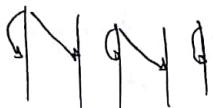
1. NEED MULTIPLE ACCEPTORS
but, we may run into two concurrent proposals with no causality \rightarrow total order anomaly



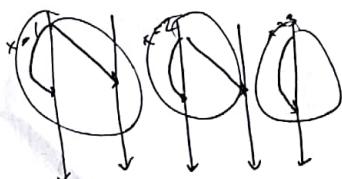
majority
to accept, we need multiple acceptors to accept in order to decide what to do.

2. Acceptors need to accept multiple values.

majority anomaly



try:
 ① ~~1 acceptor~~ \times didn't work
 ② Acceptors accept first didn't work
 after they see.

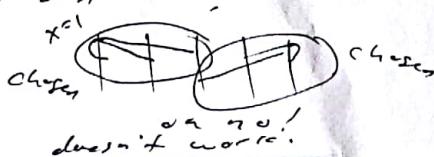


- no majority vote.
so we cannot decide.

• each node is a proposer & an acceptor.

• we can send accept to ourselves
why? acceptors sometimes need to accept multiple values.

⑤ Last value?



doesn't work!

3. NEED A TWO-PHASE PROTOCOL

4. NEED TO ORDER OUR PROPOSALS
highest ~~extreme~~ numbered gets chosen.

NUMBERING PROPOSALS

1. EACH NUMBER MUST BE UNIQUE
2. MUCH INCREASE MONOTONICALLY

Import clock with identity,
 counter / id
 typically done
 by left shifting a 64-bit
 number.

sufficient
 to implement
 Paxos.

want capability
 new proposer to
 come along & unlock
 system with slow
 nodes.

PAXOS

Proposers

1. CHOOSE A NEW PROPOSAL NUMBER
2. BROADCAST A PREPARE MESSAGE CONTAINING (1) TO ALL ACCEPTORS (SIMILAR TO CMF IN 2PC)

3. WHEN A MAJORITY OF ACCEPTORS RESPOND, PROPOSER CHOOSES VALUE TO BE ACCEPTED.
 → HER OWN VALUE, IF RESPONSES ARE UNRESTRAINED.
 → THE HIGHEST NUMBERED PROPOSAL, OTHERWISE
4. BROADCAST ACCEPT(N, VAL)

5. WAIT FOR MAJORITY
If min-proposal = N
 from every response,
val is chosen.
 else GOTO 2.

Acceptors

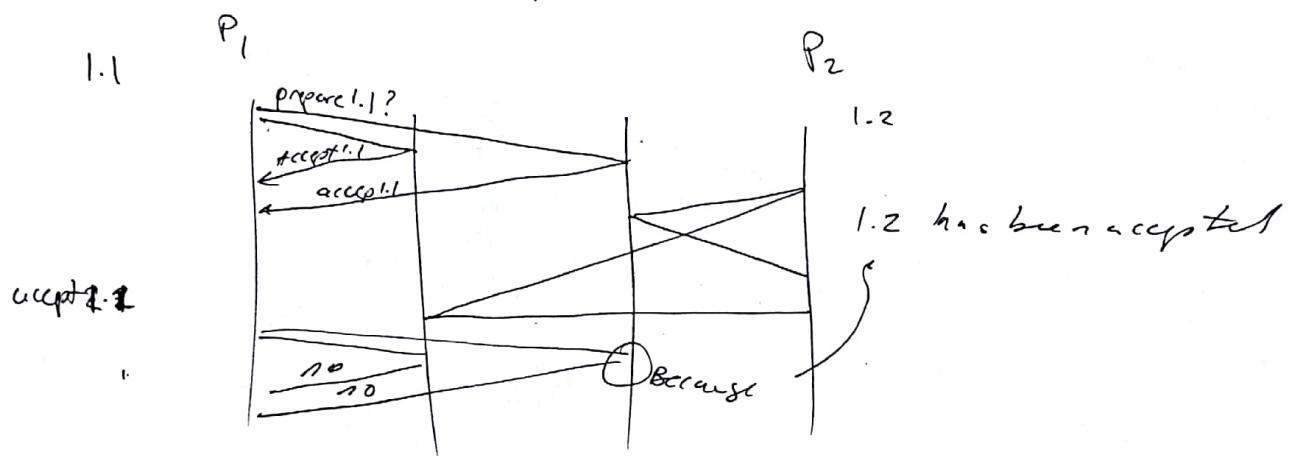
- | |
|------------------------|
| Variable: min-proposal |
| Accepted Proposals |
| Accepted Values |
3. IF N IS GREATER THAN MIN-PROPOSAL
 → SET MIN-PROPOSAL TO N.
 REGARDLESS OF CHANGE TO MIN-PROPOSAL
 RETURN (ACCEPTED PROPOSAL, ACCEPTED VALUE)
 5. ~~REJECT~~ IF $N \geq \text{min-proposal}$
 then, → accepted proposals & min-proposal
 are N.
 ACCEPTED VALUE = VALUE.
 RETURN (MIN-PROPOSAL)

Fisher Lynch Patterson Result (FLP)

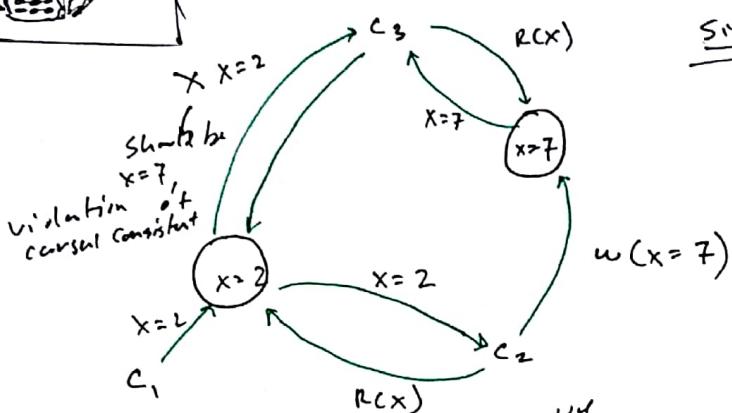
FLP IMPOSSIBILITY RESULT:

In an async system, any deterministic consensus protocol has executions that fail to terminate.

AKA: "consensus is impossible"



11/14/17



Single Variable
example

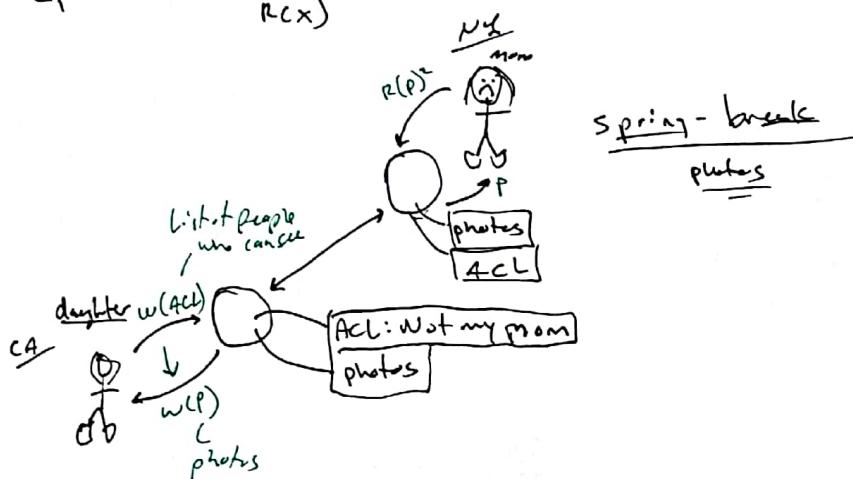
NOTE: $C_1, w(x=7)$
 $w(y=2)$
 $C_2 r(y=7)$

we know it
 C_2 reads $y=7$,
Somewhere in
time $w(y=7)$
occurred.
suppose:

$$C_2 w(z=2)$$

$$C_2 \neg w(z=2)$$

if $C_2 \vee G(z=2)$
we can see
all events which
happened before.



Spring-break
plutos

<u>MIDTERM BD</u>	
Median:	70.5
mode:	68.5
mean:	78
STDEV:	15.6

PAXOS

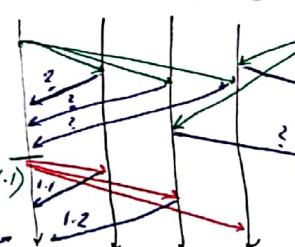
prepare(1.1)

① will
try again if
no majority.

accept(1.1)
($x=1$)

oh no!
need to try
again.

P. A₁ A₂ A₃ P₂



Prepare(1.2)

goal: No matter the order

Accepter nodes get writes,
at the end, they're all
friends who agree
on the value.

accept(1.2, $x=2$)

Since we need
a majority

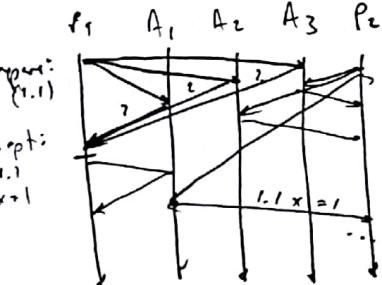
②

Propose(1.1)

accept:
1.1
 $x=1$

Prepare(1.2)

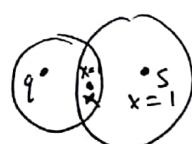
$x=0$



is over cautious, but necessary
for correctness

/ changes its
own value to
accept(1.2, $x=1$) value in
system.
its update
number

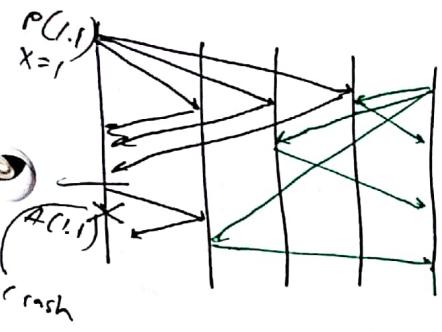
r
s



$x=1$ is chosen due to
two majorities may overlap
on a single node.

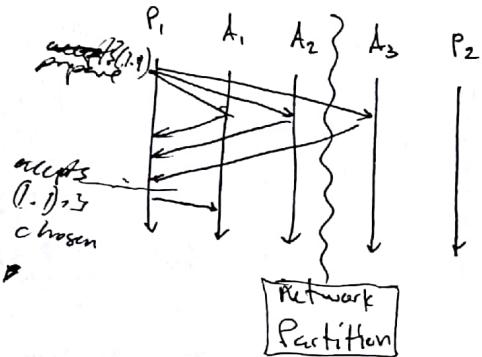
• must assume

if we see a single
constrained value, we
use that value as a
Safety inf. choose higher
number to satisfy Liveness



$\text{accept}(1.2, x=0)$
due to fact that
A crashes.

network partition example:

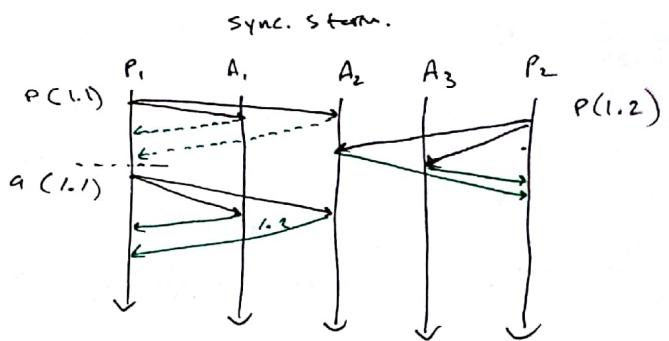


so, we can implement Paxos

Random timeout
exponential backoff
Leader election

LEADER ELECTION

1. ALL NODES PERIODICALLY BROADCAST THEIR ID TO ALL NODES
→ EVERY T ms
2. IF I HAVEN'T HEARD FROM A HIGHER ID (IN S ms)
→ Act as a proposer
3. IF I HEAR FROM SOMEONE WITH A HIGHER ID,
→ Act as an Acceptor.



- John Osterhout - Stanford Paxos
 - need to re-read paxos notes
 - simple.
 - Need to know how to run paxos by hand for final!
- FINAL**

P_1 finishes Propose phase
but P_2 finishes Propose phase
before P_1 accept phase completes,
and we're stuck in a bind.

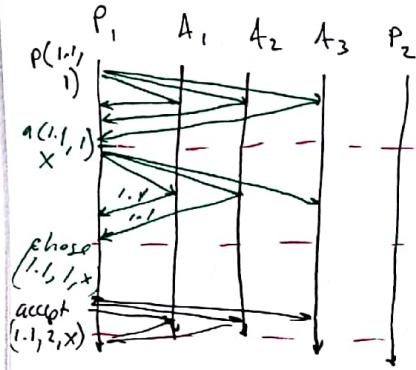
→ To solve: introduce
timeout to solve issue of
this loop. timeout is randomized
gives other proposer chance
to finish, this allows me
to do what I want. We
need to choose random
 timeouts to break out of
Synchronization storm.

Pick a number 0-10ms
Sleep, if we get blocked
again → 0-100ms, every
block, we increase exponentially
to make condition ~~more~~ unlikely
exponential backoff

(3)

Prepares are there to block earlier proposals in which they shouldn't accept.

TERMINATION
MESSAGE Complexity
Config changes
But wait!!!



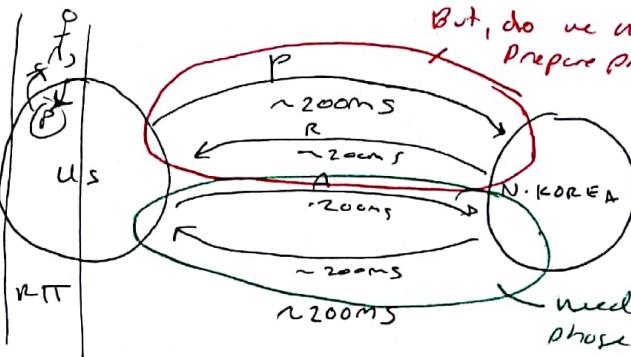
ONLY NEED to Prepare one time each time we're the leader.

Steps down if a higher round number is seen by us.

→ In which case, he drops back to prepare phase & chooses an extreme.

2 WAN RTT

MULTIPAXOS



P₁ A₁ A₂ A₅ P₂

note: delays near 250ms will be visible to user.

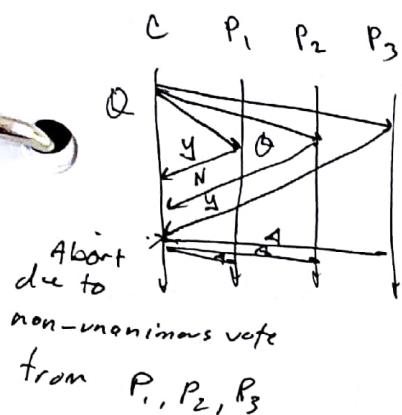
WIDE AREA NETWORK ROUND TRIP (\geq TIMES).

- need accept responses to ensure outcome.
- REDUCE MESSAGE complexity from 2 RTT to 1 RTT.
- WE CAN AVOID THE PREPARE PHASE AND GET AWAY FROM IT.

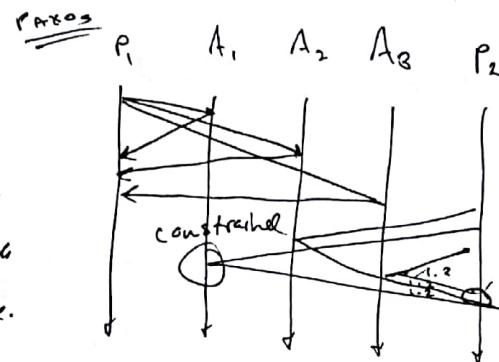
2 RTTs
for worst case
1 RTT
for normal/average case

• Pass round number & slot number
Prepare (1.1, slot 1) Should be called/referred to as phase round. which round of prepare

2PC + CP

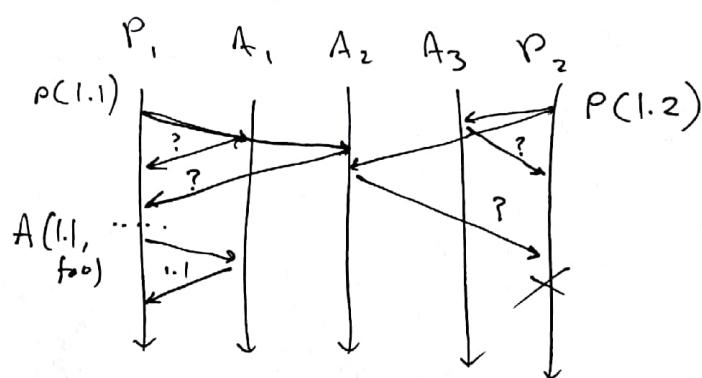


we know, this is a case where CP can decide Abort to non unanimous vote.



A_i:
Min proposal accepted P accepted V

Since Consensus for task crashes,



①
WEAK CONSISTENCY
• SESSION GUARANTEES

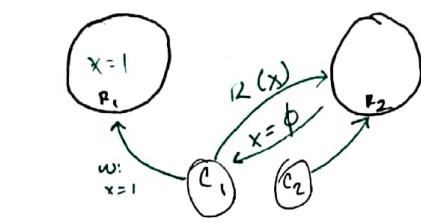
each Replica
needs to maintain
a view of the world.

if a client writes
to a node & reads
from another.



what if we read from
a node & it doesn't
have that key?

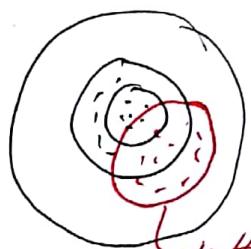
block or
forward.



Anomaly with respect to
"RETD your writes"

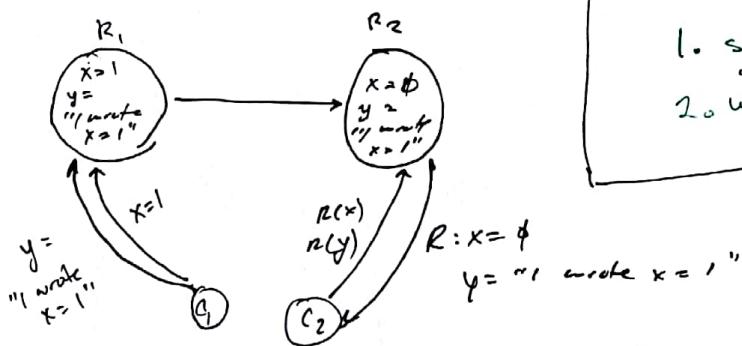
"what could possibly go
wrong?"

→ Ruling Out Anomalies



what could go wrong

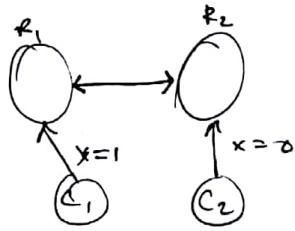
1. STALE READS
written & reads conflict
2. WRITE CONFLICTS



could use stickiness
client → R_1 for session.

Easy Trick:

KEEP VALUE ON CLIENT.
CACHE WRITES
PERFORMS WRITE.



* what if
we have
an echoing
reliable
broadcast?

1. READ YOUR WRITES.



2. MONOTONIC READS

→ CAN'T READ OLDER VALUE

3. WRITES FOLLOW READS

- RESPECTS "HAPPENS BEFORE"

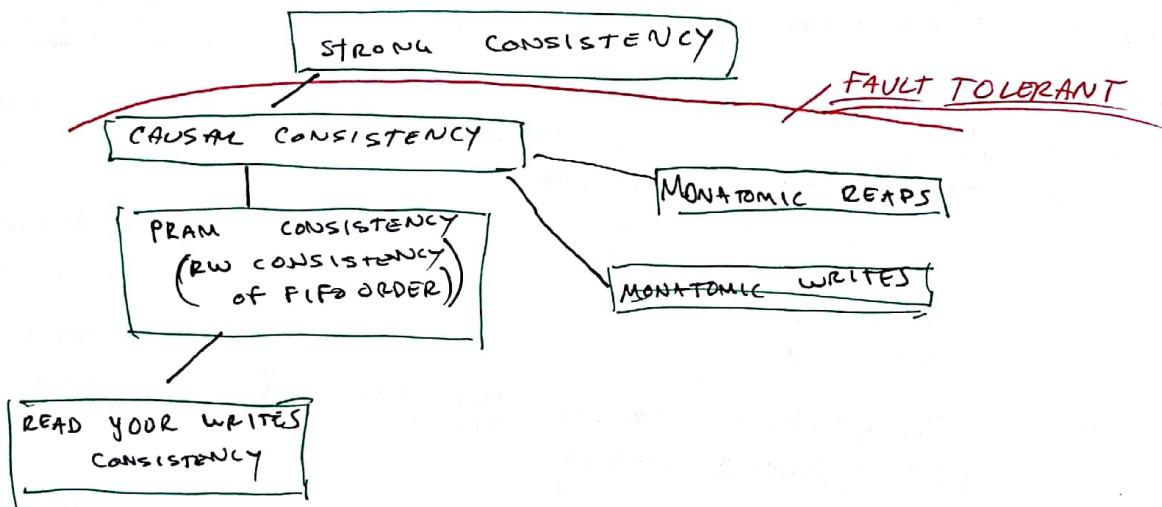


4. MONOTONIC WRITES → "WHO COMES LAST?" - P. ALEXANDER

SESSION GUARANTEES

NOTE but doesn't
our Asy3 use
causal consistency
yet still allow
stale reads?

"Cache
Rules
Everything
Around
Me!"



11/21/17

TODAY'S

TOPICS: PERFORMANCE

TRANSACTIONS &

Concurrency Control

- conflicts
- schedules
- serializability
- deadlocks
- 2PL, OCC
- Recovery
- ACID

PARALLELISM & SCALING

- Speedup - Scaleup
- Scaleout - Amdahl's Law

PLACEMENT OF DATA & COMPUTATION

- Partitioning Schemes

Amdahl's Law:

$$T(N) = \beta + \frac{(T(I) - \beta)}{N} \quad \text{PARALLEL EXECUTION}$$

SERIAL EXECUTION

(CAN'T MAKE FASTER. THIS IS THE BOTTLENECK)

TIME TO RUN ON ONE MACHINE

RESOURCES

$$T(N) = \frac{\beta}{S} + \frac{T(I) - \beta}{P}$$

Speedup: measures how the difference in runtimes with optimized upgrades

TIME UNTIL JOB COMPLETION

"Fixed Workload"

looking at effect of optimization due to inclusion of resources.

RESOURCES

Scaleup:

describes ~~a~~ single site system — adding cores to a box.

Completion time

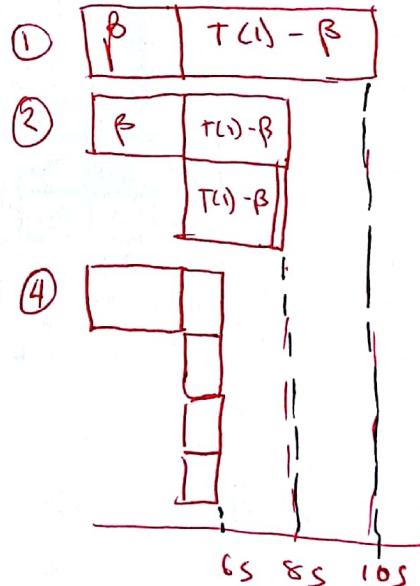
PERFECT SCALEUP

— Practice

— Theory

WORKLOAD & RESOURCES

Scaleout: increasing size of cluster
(actually adding CPU's)



BIGGEST LWS → TURN OF CONSISTENCY MECHANISMS THAT INCREASE LATENCY.

- OPTIMIZE β → SERIAL EXECUTION.

(2)

DATA PLACEMENT

- Computation → easy.
- Data → non-trivial, (map reduce, processing large amounts of DATA)

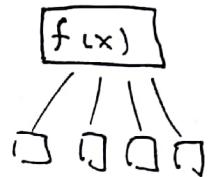
PARTITIONING:

PROBLEM OF DATA PLACEMENT

1. **UNIFORMITY**: AS WE INGEST DATA, WE DISSEMINATE DATA EVENLY. (HASHING?)

- DIVIDE DATA COMPLETELY EVENLY = QUICKEST SOLUTION → NO EASY!!

2. **SCALABILITY**: AS $\rightarrow \infty$, WHEN POSSIBLE, WE WANT LOCALIZED AS MUCH AS WE CAN.



3. **WRITE**: SUPER FAST } HARD TO DO BOTH IN PRACTICE

4. **READ**: SUPER FAST }

5. **STABILITY**:

A. RANDOMLY → NOT A BAD PLACE TO START.
↳ Kicks Ass if you never need to REAP DATA

B. HASHING → Just as much uniformity as $h(x) \mod k$ gives us bucket random & its deterministic !!

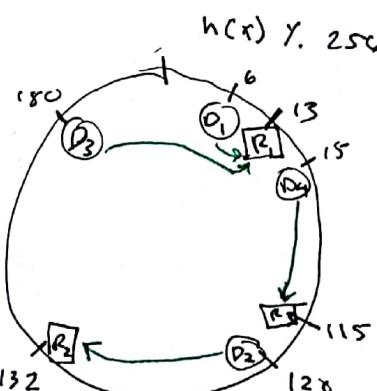
C. RANGE → LEADS ARE REALLY FAST, SUPPORTS RICHER QUERY.

D. ROUND-ROBIN → GREEDY → ALMOST STATELESS

E. DIRECTORY

F. CONSISTENT HASHING

if R_3 crashes
Data will look
for successor on
circle, now we
just need to remap
one key.



$$D_1, D_2, D_3, \dots$$

$$R_1, R_2, R_3, \dots$$

$$h(R_i) \mod 256 = 13$$

$$h(D_i) \mod 256 = 6$$

CLUSTER FILE SYSTEM DEVICES
UCSC / CEFX SERVER

Given a cluster map, we can give it to a function & find out where the data (relatively) needs to be stored.

BIP MAP INDEX OVER POINTER

③ CTPS 128

1. ESTABLISHING TOTAL ORDER

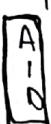
A - Atomic

C - consistent (mibble mibble)

I - Isolated ~~deadlock~~

D - Durable

NOTES: because no-one has defined consistency is dB.



T ₁	T ₂
X = R(c)	X = R(c)
X = x+1	X = x+1
W(c=x)	
(x=2)	W(c=x)
	(x=2)

Two transactions happened, but we only incremented the counter once.

NOTE: In a concurrent schedule have to think about what happens ~~when~~ between reads & writes completions.

SERIALIZABILITY → Simultaneously Stronger & weaker than total ordering.

Conflicts:

write-write conflict

2 xacts wrote to same variable

write-read conflict
read & write

T ₁	T ₂
w(x=1)	w(x=2)

SCHEDULE

AN INTERLEAVING OF PRIMITIVE OPERATIONS.

T ₁	T ₂
w(x=1)	w(x=2)
...	...
	r(x)

EQUIVALENCE

- 1 VIEW EQUIVALENCE: TWO SCHEDULES ARE VIEW EQUIVALENT IF THE RESULT IN THE SAME END STATE OF DB & ALL READS RETURN THE SAME VALUE.

NP HARD

w(x=2)	w(z=bad)	OK
r(x=2)	w(x=2)	OK

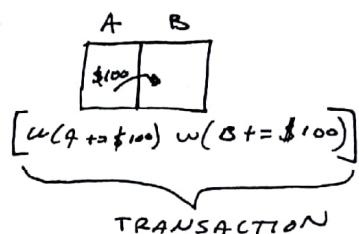
TERMS:

conflict: two xact conflict if they operate over same variable & one of them is a write operation.

here to be over same unit or row & one has to be a write

- 2 CONFLICT EQUIVALENCE - implies view equivalence if conflicting operations occur in same order but not the same in two schedules. other way around

A SERIAL SCHEDULE → ALL TRANSACTIONS RUN TO COMPLETION BEFORE THE NEXT BEGINS.



ok for none or both to happen, but not one or the other.

(4)

CMPS128

SERIALIZABILITY

operational transform
google docs
uses google
wave
-granularity

- A schedule is serializable if it is equivalent to some

Serial schedule - duh!

order in which conflicting Xacts happen is
equivalent to

Two-Phase Locking → "Strict 2PL"

all you need to do:
carry solely
about Xacts
2 place order
on them.

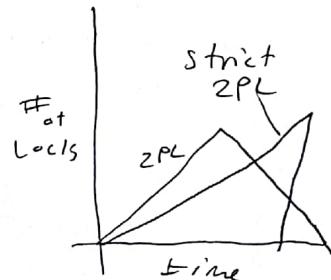
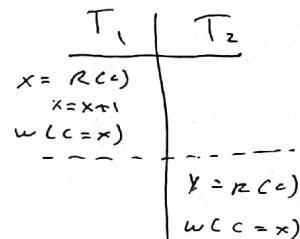
1. BEFORE writing/reading a variable, Acquire a w/r lock

2. Hold all locks until end of Xact

3. R×R - are compatible

wR - conflict

ww - Conflict

1. DEADLOCK DETECTION2. DEADLOCK REMEDIATION

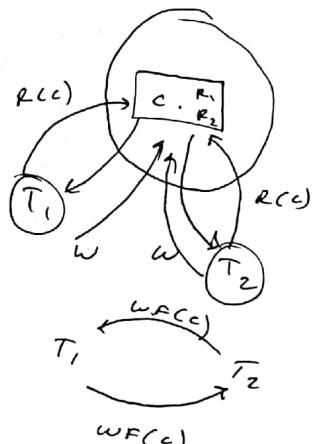
→ STOP 'em → Abort/pick a victim. — never

→ PREVENT 'em → (111)

↳ USING A TOTAL ORDER

- Xact.

↳ shouldn't
penalize older
guy.



now we're in a
deadlock.

Assumes
Infinite
Resources

"Ask for
forgiveness"

Conflicts
Abst. throwing away work.

Optimization Concurrency Control
assume conflicts are rare.

"You down with OCC?"

2PL → Ask for Permission

OCC → Ask for Forgiveness

① Someone dies &
has to restart.
② ALL TRANSACTIONS

③ RUN WITHOUT SYNCHRONIZATION

④ Xacts write to 'local' copies of variable,
On commit → Validation occurs

example {
T₁: Read-set {x} } if { R_{S1} ∩ W_{S1} = ∅ }
write-set {x} }
{ T₂: Read-set {x} } if { R_{S2} ∩ W_{S2} = ∅ }
write-set {z} } if { R_{S2} ∩ W_{S1} = ∅ }

④ commit local w, test w. concurred
all other
concurrent
Xacts.
quick lock

Distributed naming

- DNS
- DHT's & P2P
- CHORD

RECAP:

TRANSACTIONS

- Concurrency Control
 - 2PL: Ask for Permission
 - OCC: Ask for forgiveness

Serializability

- T_1 T_2
 T_3



→ MAKES SURE END STATE OF DB IS SAME

"View Serializability" - effects are indistinguishable
"Conflict Serializability" - ~~same~~,
- conflicting operations

DISTRIBUTED NAMING

why we name things?

so humans can read/tell things apart

• Identity - disambiguate

• Addressing / Retrieval / data content.

• names behave as indices.

• SERVICE DISCOVERY... (name is task)

• Indexing

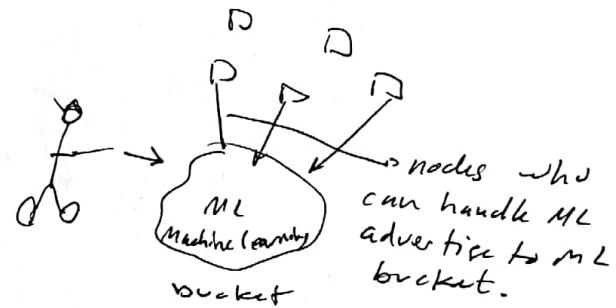
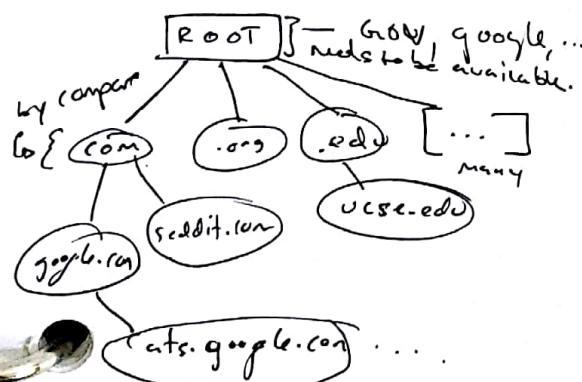
• Decoding

• machine learning

APPROACH

handles case A, B

- HIERARCHICAL NAMESPACE



DNS - Your GPU is a replica!

A SCALABLE

B Multiple Administrative Domains

C Low LATENCY (aggressively cached) (even in local op)

D AVAILABILITY is King.

• DNS is weakly consistent
- "imagine if every laptop had to ask"
- That's what a strong consistent DNS system would suck my ass!" - Peter

E Geo-Scale

11/30/17

BIG DATA

- Data Parallelism
 - make stuff go fast

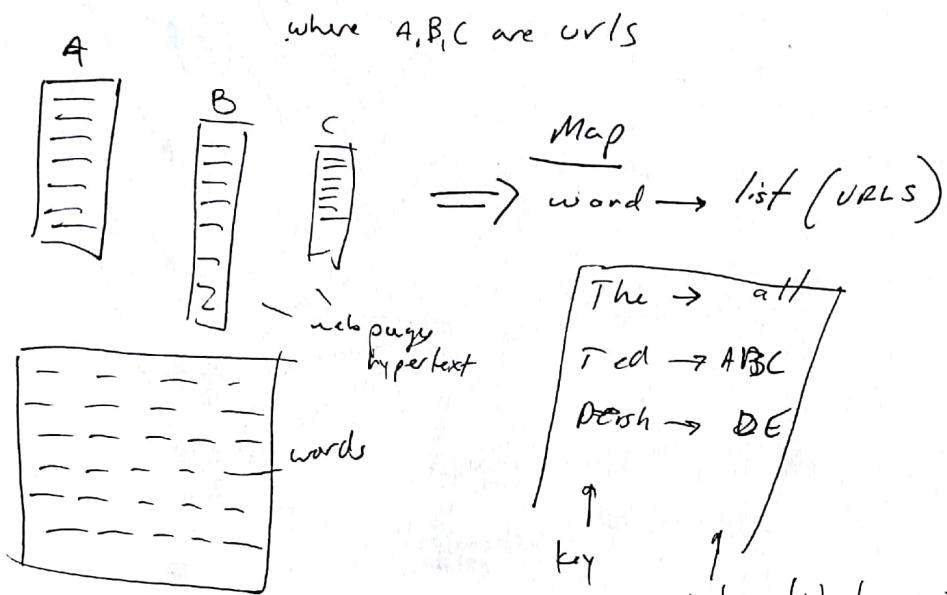
INVERTED INDEX



Searched term: Ted Dersch

$$\text{map}(\text{Ted}) \cap \text{map}(\text{Dersh})$$

$$A, B, C, D, E = A.$$



on a single machine:

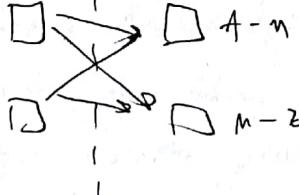
for every url

if word in url.doc

append url onto values list for word

many machines: Asg 4

all-to-all communication



mapReduce

filtering, transformation phase — map
combination — reduce

$\text{map}(k_1, v_1) \rightarrow \text{list}(k_2, v_2)$

$\text{Reduce}(k_2, \text{list}[v_2]) \rightarrow \text{list}[v_2]$

map: for each doc

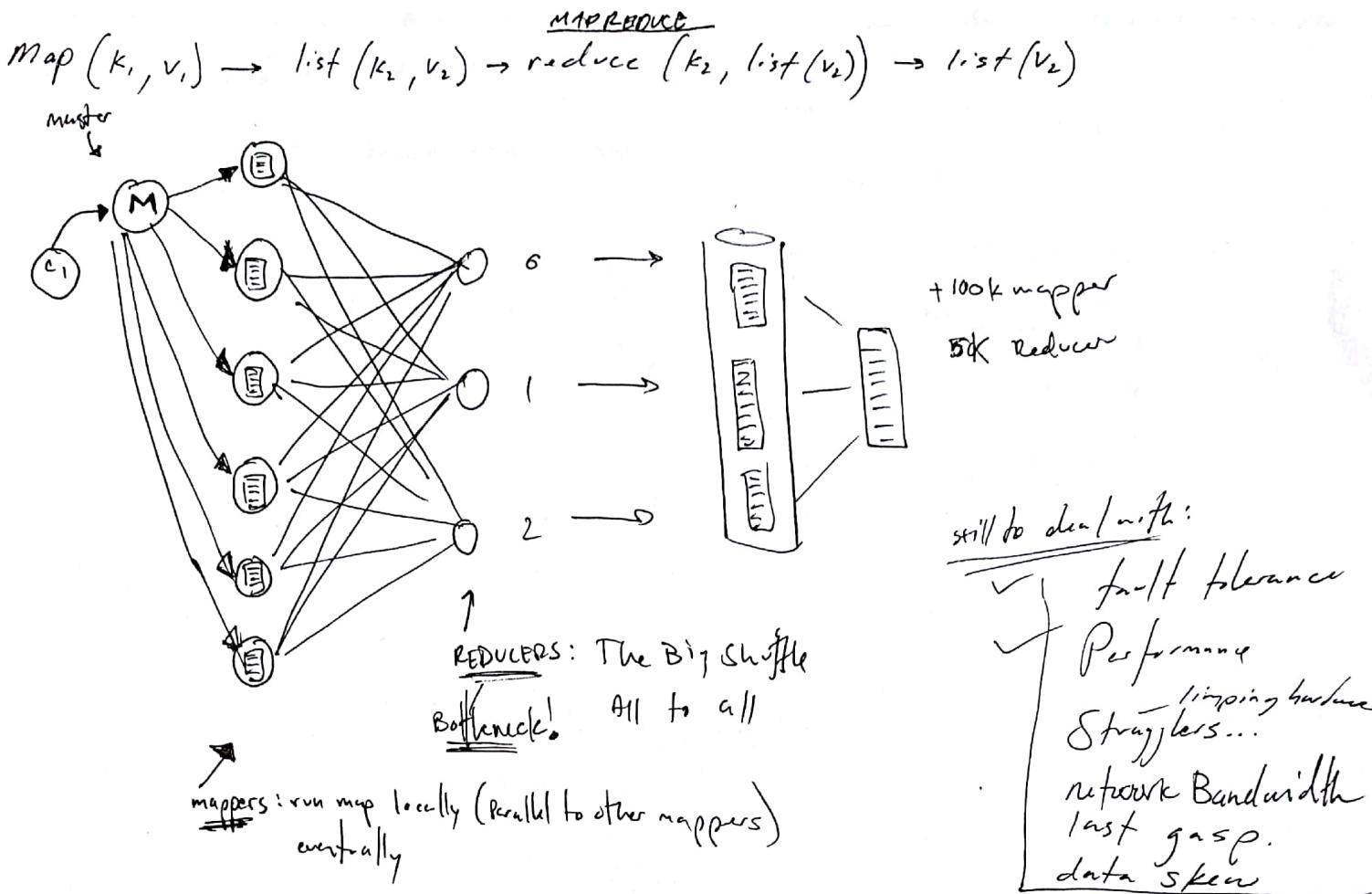
for each word

emit $\langle \text{word}, [\text{doc}] \rangle$

reduce:

for k_2 for every pair $\langle k_1, [v_1] \rangle$
 $\langle k_1, [v_2] \rangle$

emit $\langle w, [v_1 + v_2] \rangle$



☒ PERFORMANCE → Problem of Locality

- 1) SCHEDULING COMPUTATION CLOSE TO DATA
(RATHER SEND A FUNCTION THAN 1TB DATA)

☒ STRAGGLERS

Insights: limping hardware

Speculation: Predictive Branching.

if a machine is taking a long time to complete a task, tell another process to run task.
takes care of old computers/slow processes.

Bonus: gets rid of stragglers

Neg: do ~ 1.1 the work ($\times 2$) work.

TRADEOFF: throughput for latency.
"resources for time"

☒ FAULT-TOLERANCE

NOTE: mappers ack master
timeout is sufficient
mappers - promote new
it comes back
tell it to fuck off
Reducer - promote new & give
it same bucket number
since mappers store
file upon completing
(NFS lookup) to get
new reducer up to speed.
master - Boom Analytics
masters get smire with
Paxos.

③

We, databaseers, have been doing this for ages, we just didn't have a ~~fun~~ cute name for it until we didn't work at Google.

1. RESTRICTED PROGRAMMING MODEL SQL

2. DATA PARTITIONING

3. DATA SHUFFLING

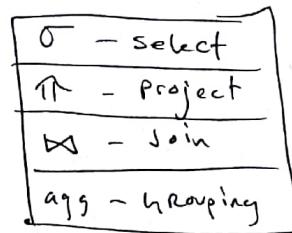
Select $x, y, \text{count}(z)$
from R, S

where $R_x = S_y$

and $R_x < 10$

group by x, y

agg ($\sigma_{x \in 0} R \bowtie S$)



SECTION

ROLESMAPPERS:

PARTITIONING

DATA
fault tolerance
TASKS
consistency

LOCALITY

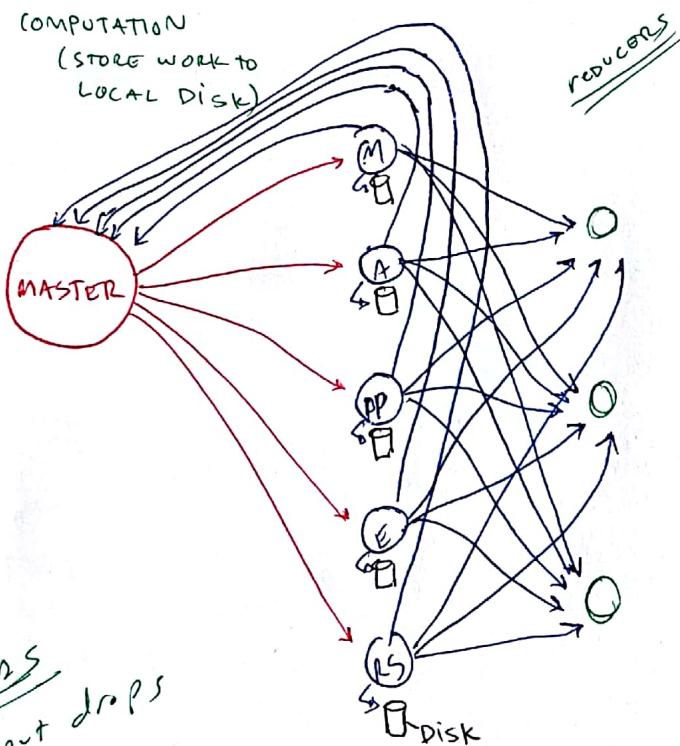
distance to data needed.

REPLICATION

OF DATA

OF COMPUTATION

(STORE WORK TO LOCAL DISK)



STRANGERS
Throughput drops due to more than one node on a single bus. + allocation of double + allocation of work tasks.
Latency drops avoid K timeouts

REDUCERS:

INSERT DATA FROM MAPPERS

AGGREGATE USING JOIN

MASTER:

(MAPPERS)

DETERMINES WHO GETS WHAT WORK

DETERMINE WHICH REDUCERS MAPPERS TALK TO.

ALL - TO - ALL COMMUNICATION

Information Disseminated By Master

- which mappers get which work?
- REDUCERS LEARN WHERE TO SOURCE THEIR DATA
- KNOWLEDGE OF WHO IS UP.

FAULT TOLERANCE

- Master is replicated, Paxos
- Mapper's detect: timeout remedy: allocate node, if the same mapping values
- Reducers: detect: timeout remedy: new reducer at same bucket & same source.