

(doesn't the digital clock work on batteries?)
↳ req power, isn't it?

classmate

Date 28/07/25

Page _____

✗ they hardware? or process? or advice? (combination of both)

i) Distributed Systems: hardware components as networked computers

which communicate by message passing

→ collection of independent computers (hardware) or processes (soft)

→ appears to be coherent system (wire) → as to achieve a goal

→ multiple # computers, interconnected & shared system (3 features)

→ failure prone: as all of them try to achieve the same goal

→ it is layered over networks ⇒ for each layer specific protocols present //

ii) Perfect synchronization not possible ← due to delays in networks

→ one obj. per TCP
non-persistent connection
persistent "

autonomous, programmable, asynchronous

iii) HTTP

& failure-prone

→ Stateless (as complex if they have to manage) - servers maintain no info

→ uses TCP transport service about past client requests

iv) Issues Related to Distributed Systems

correct

• No global clock: no single notion of time (synchronization)

- unpredictable failure of components

- highly variable bandwidth

- Possibly large & variable latency → (in Internet, it's high)

→ for memory (seconds), it's in ms

• large no. of hosts challenges concurrent execution, independent failures*

Autonomous administration, heterogeneous

v) Network communication → Challenges → A synchronous, unreliable, incomplete

Common goal - (consistency)*, transparency

* real issues

vi) Middleware - (diff software b/w) (why by DS req it?)

• Software programs

→ different functions exist

• Why necessary? - Distributed systems: functional distribution

inherent

→ application differs

VII) Challenges & Goals of Distributed Systems

→ ex. remote method invocation, group communication, replication

- Heterogeneity — can be handled by Middleware

↙ a SW layer providing a programming

(variety &

abstraction as well as masking the heterogeneity

difference of

of the underlying networks, HW, OS, & programming languages

underlying networks

infrastructure)

Software layers
in DS



- Openness — determining whether system can be extended & implemented in various ways

- Security — DS must protect shared info & resources

- Scalable — Can handle additional no. of users/resources without suffering noticeable loss of performance

- Failure Handling — Failure in DS is partial

(try Redundancy) ↙

↑

making handling of failures difficult { some components fail while others continue to function

How to handle? →

(i) Masking failure — trying retransmission

(ii) Defeating

(iii) Tolerating

How to recover? → (i) Rollback

(ii) Undo/redos in transmission

- Concurrency — multiple users — stale data in shared system

↳ can lead to deadlock.

- ~~Hetero~~ Transparency - concealing the heterogeneous & distributed challenges? | nature of the system so that it appears to the user like one system

Categories

(i) Access - access local & remote resources using identical operations) (....)

28/07/25

Topic : System Model \rightarrow description

(User roles, hardware \leftarrow (In high level)
hardware, & have responsibilities divided) \star

i) Architectural Model : Structure of System in terms of components

mapping of processor to processor
server & client not 1 to 1 as a server can provide many services

tells what

Distributed System Architecture - 2 types

Client - Server Architecture

V/s Distributed Object Architecture

\rightarrow distinction b/w Client &

\rightarrow no distinction

servers (note: there are not of different machines)

machines)

any one can ask /
provide

\rightarrow tells how the system behaves

Fundamental Models - a system model

i) the entities

: defines

\leftarrow

Purposes: (i) explicitly states all

ii) their interactions

relevant assumptions

iii) characteristics that affect their

(ii) generalise what is possible or

behavior

not based on the assumptions.

* Concerned with properties common to all architectural model

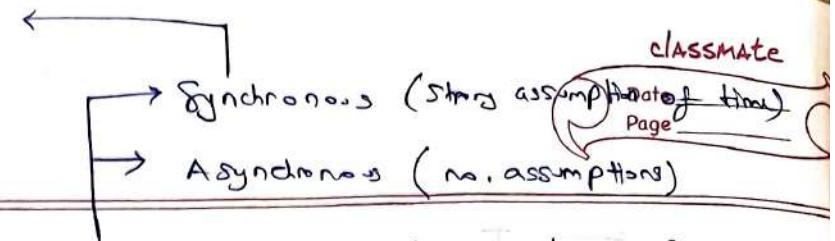
* often slides for
3 key features

{ (i) No global timer in a DS

(ii) communication occur by messages \rightarrow vulnerable to attacks, delays etc

(iii) failure prone

bounds on processing time (min & max), msg transmission delay, local clock drift rates



- 3 models are : (i) Interaction model - deals with performance & the difficulty of setting time limits in a DS

types of failure

(i) omission

(ii) arbitrary

(iii) timing

(ii) Failure model - attempts to give a precise specification of the faults that can be exhibited by a process & communication channels
(message eventually reaches the destination) : validity \leftrightarrow integrity : no manipulation of message
(iii) Security model - discusses the possible threats to processes &

Communication Channels

achieved by :

(i) access rights

(ii) principal (??)

introduces the concept of a secure channel, which is secure against these threats

- Performance of communication channel maintained by :
- latency : sender to receiver (time taken) in sending

- includes propagation delays - propagation in a network
network delays - delay in accessing the network
OS delays - time taken by OS

- bandwidth : capacity of communication channel

- jitter : variation in time taken to deliver a series of messages

Failure Types

can be marked

(i) Omission failure

Crash stop (fail-stop)

process

Crash recovery

communication

Send omission : loss before

Sending process & outgoing msg

buffers

Channel omission

receive !!

(ii) Arbitrary failures

Byzantine

failure

caused by NITM, viruses, worms etc.

process - omits intended processing steps

or takes unintended processing steps

Channel - messages corrupted, duplicated,

large delays etc.

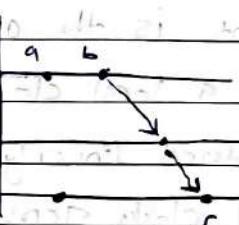
Topic : Time in Distributed System

- Berkey Algo : server after computation sends the offset & not the time to its clients
 - NTP Network Time Protocol : Best works with Symmetric mode
 - causality : means cause & effect
 - 2 local events in 2 different machines can't be directly compared unless explicit relationship
- and also what about solar shadow also?*

• Physical Clocks

i) International Atomic Time (TAI):

no. of ticks of cesium 133 atom since



relationship b/w a & c?

1/1/1900 (atomic second) \rightarrow Can't say!

ii) Atom Clock : One second defined as (since 1967) 9,192,631,770 transitions of the atom Cesium 133

iii) Slowdown of earth $\xrightarrow{\text{leads to}}$ introducing leap seconds \Rightarrow Correction of TAI is called Universal Coordinated Time (UTC)iv) if $C_i(t)$ be a perfect clock(i) Clock $C_i(t)$ is called correct at time t if $C_i(t) = C(t)$ (ii) $C_i(t) \approx C(t)$ if $\frac{d(C_i(t))}{dt} = \frac{d(C(t))}{dt}$ (iii) 2 clocks $C_i(t)$ & $C_k(t)$ considered synchronized at time t if $C_i(t) = C_k(t)$

• Clocks

 \rightarrow computers contain physical clocks (crystal oscillator) \rightarrow Clock skew : Instantaneous difference b/w the readings of two clocks~~→ Clock drift~~ : Crystal based clocks count time at different rates?

↳ divergence

 \rightarrow relative difference in clock frequencies (rate) of 2 processes

- Why Synchronization? — correctness
fairness

- Process: present in a synchronous DS

- + has a state
- + takes actions which change its state
- + Event is the occurrence of an action
- + has a local clock \Rightarrow events within a process can be ordered linearly
- + but clocks across processes (not synchronized) is not

- Max drift rate b/w 2 clocks with similar MDR is $2 * MDR$ (MDR)

NOTE: A non-zero clock drift causes skew to increase (eventually)

- Synchronization \rightarrow External (synchronized with an external source of time) or Internal (synchronized with one another)

$$\frac{dc}{dt} > 1 \quad \text{clock ticks faster than UTC}$$

$$\frac{dc}{dt} = 1 \quad \text{clock ticks at UTC}$$

$$\frac{dc}{dt} < 1 \quad \text{clock ticks slower than UTC}$$

(Check formula part again)

- for synchronization (Asynchronous DS) \rightarrow Cristian's Algo
- + Berkeley Algo
- + must be with some delay for internal synchronization

NTP (Network Time Protocol)

if 2 processes needn't req exact time

- Logical Clocks - kind of says ordering more important
 - ↳ provides timestamps //
 - but if clocks need to have sometime \Rightarrow physical clock
 - \rightarrow was the happens before relation
 - \rightarrow assigns logical timestamps to events

Lamport Logical Clock - each process has a local counter

- (on each event it increments its clock)
- When sending a message, it sends its clock too
- When receiving, it sets its clock to $\max(\text{local}, \text{received}) + 1$

Topic: Vector Clocks, Global State

04/08/25

- Cut consistency: \Rightarrow if receive message is recorded, so sent message should also be recorded

NOTE: if $a \rightarrow b$ (events in same process) $\Rightarrow L(a) < L(b)$

but if $L(a) < L(b) \not\Rightarrow a \rightarrow b$

because a, b could be events in 2 different processes

\Leftrightarrow still be concurrent though logical clock not!!

• Vector Logical Clocks

i) address the issue as described above (\uparrow)

ii) instead of using a single logical clock, each process has N logical clocks for the other ($N-1$) processes

(kind of vector clock)

\rightarrow without synchronization, captured using causality

• Global state = Status of all processes + status of all communication channels

If & only if

it corresponds to consistent channel (i.e. messages in transit on the channels)

Topic : Chandy-Lamport Algorithm (CLA)

- Note: Channel marked as empty \Rightarrow no message has been received in there
- Linearization may not be unique (multiple possible linearizations)
- To tell happens before relationship \Rightarrow (vector comparison can be done)
- CLA results in a consistent global state
 - \hookrightarrow i.e. if e_j is included in c_{ij} , $e_j \rightarrow e_i \rightarrow e_j \Rightarrow e_i \in c_{ij}$

Prove by: fifo ordering of channels

\Rightarrow if $e_j \rightarrow p_j$ records its state
but we say $e_j \leftarrow p_j$ records its state

Since p_j recorded its state

it will send marker for first message sent to j in C_{ij} channel
is marker then message for $e_j \leftarrow p_j$ occur

$\Rightarrow e_j \leftarrow p_j$ recorded its state

contradiction

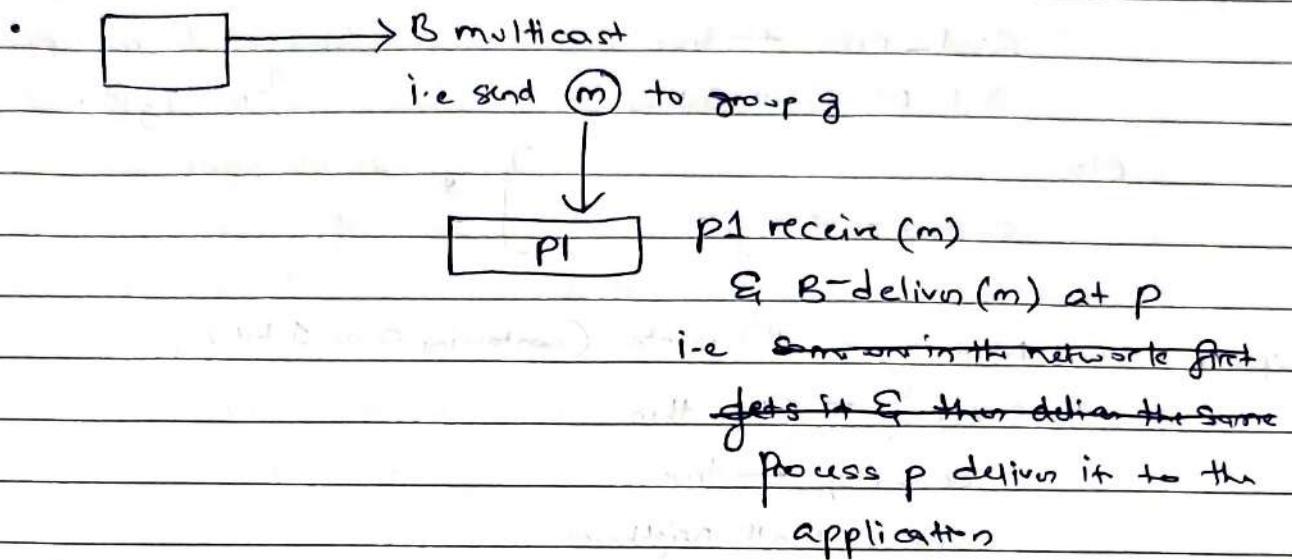
06/08/23

Previous: Synchronization // implement as well as monitoring

Topic: Brown's Communication // distributed consensus

- how to implement multiparty in DS: reliable or ordered
 - from one process to a group of processes
 - capture state + monitor + evaluate
- Distributed debugging
 - global state \rightarrow definitely
 - ? i.e. predicate \rightarrow possibly
- it is about capturing consistent global states to check whether meaningful conditions (like $|n_i - n_j| \leq 8$) actually hold during execution not just due to time artifacts.

11/08/25



- Reliable multicast

↳ Integrity: One-time sending

↳ Validity: A process which sends m correctly also delivers guarantees delivery m (+ itself) to the sender

↳ Agreement: If some one correct process delivers message m then all other correct processes in group(m) will eventually deliver m /
either all or nothing

- Order of multicast \rightarrow if multic平 (g, m) first issued

per sender \leftarrow \hookrightarrow FIFO: (m) delivered first then m' get delivered

↳ Causal: multicast $(p, m) \rightarrow$ multicast (g, m')

2 people \leftarrow $\Rightarrow m$ delivered first then m'

↳ Total: if m first then m' (for all)
for all it's same the order

FIFO (per sender)
(nothing about different senders)
Messages from the same sender are delivered in the order they were sent

Causal ordering
If one message could have influenced another, deliver them in causal order

Total ordering
All processes deliver all messages in the exact same order, no matter what!

Flooding

Initially do : If $pid = root$, then

Case 1: for 1 msg

Send msg \leftarrow true

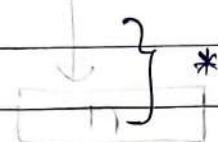
Case 2: for multiple msg

Send M to all neighbours

modify (1)

else

send-msg \leftarrow false



Upon receiving M do : it's counter (containing 0 or 1 Lit)

If send-msg = false then

send-msg \leftarrow true

Send M to all neighbours

time to complete

Time Complexity: Diameter (farthest distance b/w any 2 nodes) - O(n)

Space complexity: 2nd degree among n^2 possibilities

Message II : (Edges - O(E))

how many messages were transferred

number of messages transferred to the node

Case 3: Setting the parent pointer

A \rightarrow B

\uparrow sets parent pointer to A

↳ how many edges to point to

red modification in \Rightarrow Case 1

pid = root, then

Initially do : if $pid = 1$: parent \leftarrow root

(case 1) * \rightarrow else Send M to all neighbours

Send msg' parent $\leftarrow 1$

Upon receiving M do :

If parent == 1 then

if message 11A

parent $\leftarrow p$

(otherwise) \Rightarrow 11B

if message 11B

Send M to all neighbours

attain maximum no. of edges

and then if

root with bottom

attain shortest path also

otherwise

return (msg) n

else

To know all the children

initially do :

$$\text{non-children} = \emptyset$$

If $\text{pid} = \text{root}$, then

$\text{parent} \leftarrow \text{root}$

$\text{children} \leftarrow \{\text{root}\}$

Send M to all neighbours

else

$\text{parent} \leftarrow \perp$

$\text{children} \leftarrow \emptyset$

Upon receiving M from p do :

If $\text{parent} = \perp$ then

$\text{parent} \leftarrow p$

Send $\langle \text{ack} \rangle$ to p // to mark as child

Send M to all neighbours

else

Send $\langle \text{node} \rangle$ to p // to mark as non-child

Upon receiving $\langle \text{ack} \rangle$ from p do :

$\text{children} \leftarrow \text{children} \cup \{p\}$

Upon receiving $\langle \text{node} \rangle$ from p do :

$\text{non-children} \leftarrow \text{non-children} \cup \{p\}$

Corrections : Since everyone sends the message so eventually everyone gets it & set the children / non-children in $O(\Delta)$

Hamilton Path : path connecting all nodes visiting only once (NP complete)

Euler Path : edges visited only once

(Polynomial)

Construction of spanning tree (Algo)

 \rightarrow via BFS or DFS \rightarrow Using DFSCase 1: designated root (P_i)
(discussed)

Case 2: multiple initiators

- Initially assume: parent $\leftarrow 1$,

children = \emptyset (Only 1 Spanning tree
will
should be constructed)Unexplored includes all the neighbours of P_i (by some coordinated action)for $i \in S$ upon receiving no messageno. + 1 if $i = x$ and parent = 1 thenparent $\leftarrow i$ let P_j be a precursor in unexploredremove P_j from unexploredSend M to P_j

- Upon receiving M from neighbour P_j

if parent = 1 then

parent $\leftarrow j$ remove P_j from unexplored // as it received it fromif unexplored $\neq \emptyset$ thenlet P_k be a precursor in unexploredremove P_k from unexploredSend M to P_k

else

Send <parent> to parent

else:

Send <parent> to P_j

• Upon receiving $\langle \text{parent} \rangle$ or $\langle \text{reject} \rangle$ from neighbour P_j :

if received $\langle \text{parent} \rangle$ then

add j to children

if unexplored = \emptyset then

if parent $\neq i$ then send $\langle \text{parent} \rangle$ to parent

otherwise terminate \leftarrow for current process P_i

else

let P_k be a process in unexplored

remove P_k from explored

Send M to P_k

How do we know overall termination happened?

→ central Node which collects all terminated detail

No. of messages passed = $O(E) \approx O(4M)$

$4*E$ (worst case)

Time complexity = $O(E) \approx O(M)$

13/08/25

Constructing DFS spanning tree

w/o specified root

copy of $\langle \text{root} \rangle$ (wrt P_i) (below steps)

• Initially: $\text{parent} \leftarrow 1$, $\text{leader} \leftarrow 0$

$\text{children} \leftarrow \emptyset$

$\text{unexplored} \leftarrow \text{all neighbours}$

• Upon receiving no message:

if $\text{parent} = 1$, then

$\text{leader} \leftarrow \text{id}$, $\text{parent} \leftarrow i$

let P_j be a process in unexplored

remove P_j from unexplored

Send $\langle \text{leader} \rangle$ to P_j

(leader variable)

new-id

- Upon receiving $\langle \text{parent} \rangle$ from P_j

if leader < new-id, then

leader \leftarrow new-id, parent $\leftarrow j$

unexplored \leftarrow all neighbours of P_j except P_i

if unexplored $\neq \emptyset$ then

let P_k be a process in unexplored

remove P_k from unexplored

Send $\langle \text{leader} \rangle$ to P_k

else

Send $\langle \text{parent} \rangle$ to parent

else if leader = new-id, then

Send $\langle \text{already} \rangle$ to P_j

// otherwise leader > new-id and the

DFS for new-id stalled

(new-id) or $\langle \text{already} \rangle$

- Upon receiving $\langle \text{parent} \rangle$ from P_j :
- if received parent then:
- add j to all children

if unexplored = \emptyset then

if parent $\neq i$, then send

$\langle \text{parent} \rangle$ to parent

else terminate as root of the spanning tree

else let P_k be a process in unexplored

remove P_k from unexplored

Send $\langle \text{leader} \rangle$ to P_k

If received $\langle \text{already} \rangle$ then

{nothing to do}

Time complexity = $O(n \times m)$

↑ per process

Message complexity = $O(n \times m)$

↑ no. of nodes

Topic: Mutual Exclusion

↓ a piece of code is at most 1 client executing it at any point of time

ensured through message passing based protocols

→ entry()

→ access resource()

→ exit()

• Requirements

→ Safety: only 1 process

→ liveness: should be actively doing

→ ordering: req. granted in the order they are made

• DS mutual exclusion performance evaluation criteria

→ Bandwidth: # messages sent in each entry & exit operation

→ Client delay: (req. operation requested & granted) time duration

→ Synchronization delay: (entry & exit of process)

• In Token Ring Approach

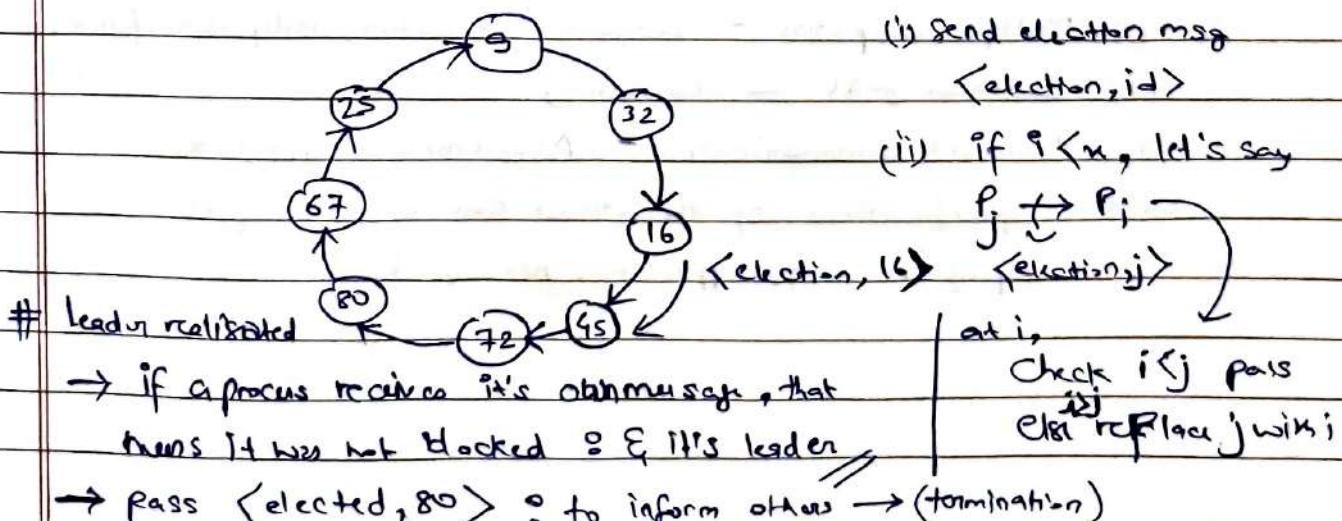
→ there is no separate entry() message

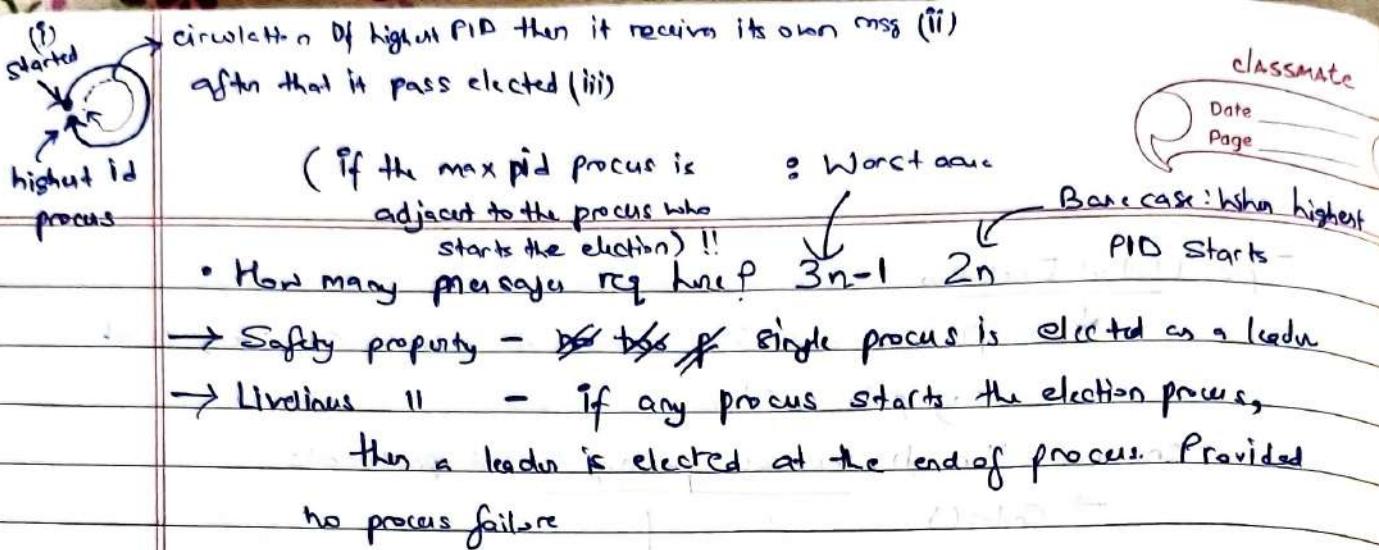
→ only has the exit() message which corresponds to token being passed

18/08/25

• # Ring Based Algorithm: for leader election

→ each process has an id





- What if a process in the ring has failed?
- Some process should detect it

MIDSEM

- L1
 - DS - hardware components / software - located at networked computers - communicating by message passing
 - Characteristics
 - (i) multiple computers
 - (ii) interconnections
 - (iii) shared state
 - HTTP Status - Server maintains info about past client req
 - ↳ persistent vs non-persistent
 - Type of distribution
 - (i) Functional
 - (ii) Inherent
 - DS issues
 - (i) no global clock
 - (ii) unpredictable failure of components
 - (iii) highly variable bandwidth
 - (iv) possibly large & variable latency
 - (v) Large no. of hosts
 - Challenging Properties
 - (i) Multiple computers - concurrent execution, independent failure
 - (ii) Common goal - consistency
 - (iii) Network communication - Asynchronous, unreliable
 - RFC 1 - specifications of the Internet communication protocol
 - specifications for applications run over them

- Heterogeneity - ensured by Middleware (cross communication & resource sharing)
- Openness - by RFCs, CORBA
- Scalability - 3 metrics : (i) no. of users/resources
 (ii) distance between the furthest nodes in the system (network radius)
 (iii) no. of organisations exerting control over the pieces of the system

Transparency - (app def) heterogeneous nature of system should be hidden

L2

- DS Architecture Type : (i) Client-Server
 (ii) peer to peer
- Fundamental Models : (i) Interaction model → Synchronous DS
 (ii) Failure " → Asynchronous DS
 (iii) Security "
- Performance : (i) Latency - delay between sending of a message by one process & its receipt by another
 includes, propagation - time taken for first string of bits to reach destination
 Network - delay in crossing the network
 Bus - t-t by OS communication services at both sending & receiving process
 (ii) Bandwidth : capacity of the communication channel
 (iii) Jitter : variation in time taken to deliver a series of messages
- Synchronous - processing time has a bound !!
 same with message transmission delay (+ local clock drift rate) → en process crash, message drop
- Types of failure - (i) omission - when a process or a channel fails to perform actions that it is supposed to do
 (ii) Arbitrary (Byzantine) Failure - any type of error
 (iii) Timing Failure

- Omission Failure - (i) send omission : b/w sending process & outgoing msg
 - (ii) channel II : loss in the communication channel
 - (iii) receive II : loss in incoming msg buffer & receiving process
- Reliable comm - defined by: (i) validity : eventual reaching
 - (ii) integrity : (msg received == msg sent)

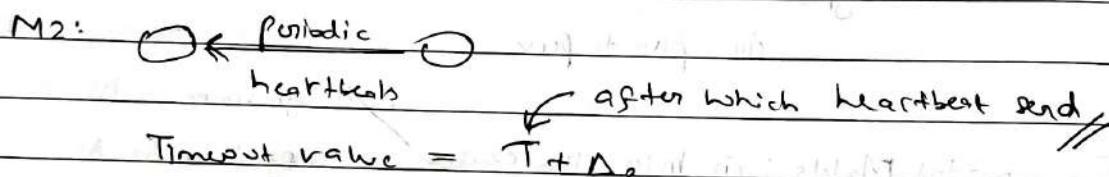
Detecting a crashed process

(i) synchronous :

$$\text{timeout } (\Delta_1) = 2 \text{ (max network delay)}$$

(ii) asynchronous :

$$\text{timeout } (\Delta_2) = k * (\text{max observed RTT})$$



Synchro - $\Delta_2 = \text{max network delay} - \text{min network delay}$
nos

Asynchronous : $\Delta_2 = k(\text{observed delay})$

Correctness of failure detection

both completeness

achievable

{ (i) completeness : every failed process is eventually detected
(ii) Accuracy : every detected failure corresponds to a crashed process (no mistakes)
in synchronous

but not in asynchronous

→ piggyback heartbeat : complete

→ for accuracy & incomplete

: never report failure ??

Worst case for failure detection

• ping ack : $T + \Delta_1 - \Delta$: Δ is t.t for last ping
from p to reach q

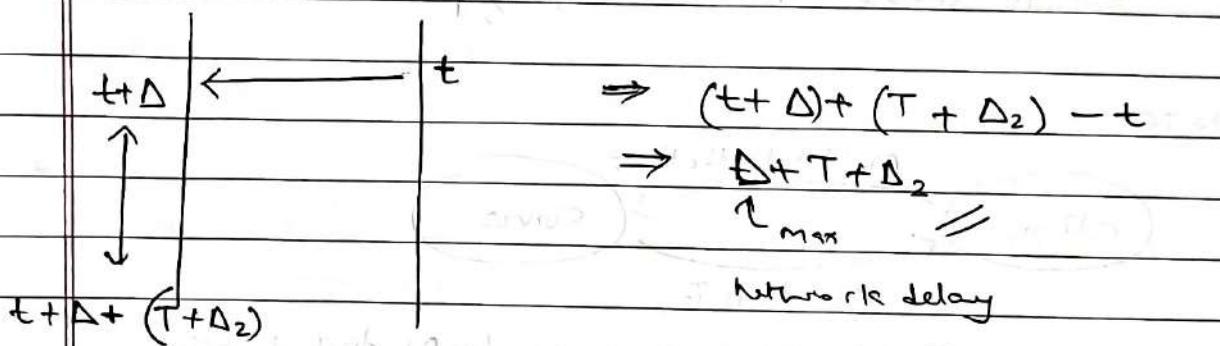
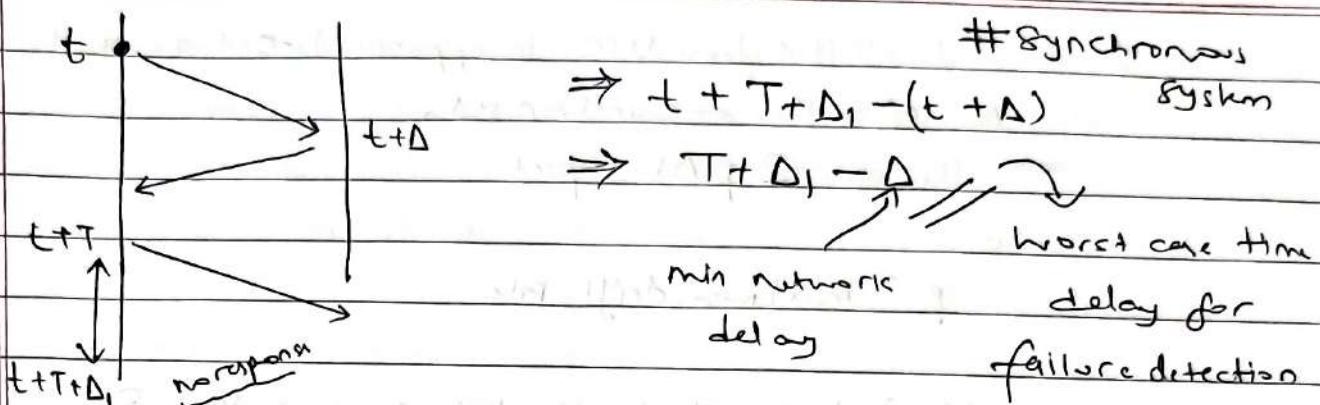
• heartbeat : $\Delta + T + \Delta_2$: Δ is t.t for last ping
from q to reach p

$$\begin{aligned} |A-B| &\leq D \\ |C-B| &\leq D \end{aligned} \Rightarrow |A-C| = |(A-B) + (B-C)| \leq 2D$$

classmate

Date _____

Page _____



NOTE: in case of asynchronous system, increasing Δ_1 or $\Delta_2 \Rightarrow$ increase in accuracy of but increase in failure time

Failure Detection

L-3

- Clock skew: relative difference between two clock values
- Clock drift rate: change in skew from a perfect reference clock per unit time
- 2 types of synchronization: (i) external : with an external clock
(ii) internal

NOTE:

$$\left. \begin{aligned} |A-b| &\leq D \\ |C-b| &\leq D \end{aligned} \right\} \Rightarrow |A-C| = |(A-b) + (b-C)| \leq 2D$$

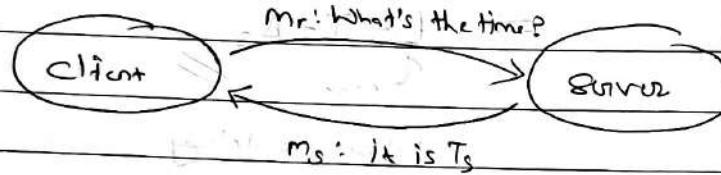
- 2 clocks if drifting from UTC in opposite direction at a time Δt after they are synchronised,
 \Rightarrow they are $2 p \Delta t$ apart

where,

p : maximum drift rate

- Also, if 2 clocks shouldn't differ by more than δ ,
 \Rightarrow it should be reset every $\frac{\delta}{p}$ seconds atleast

NOTE:



What then T_c should adjust its local clock to after receiving ms?

$$T_c = T_s + \frac{(\text{min} + \text{max})}{2}$$

$$\text{Skew}(\text{client}, \text{server}) \leq \frac{(\text{max} - \text{min})}{2}$$

How to Synchronise,

(i) Christian Algo

$$T_c = T_s + \left(\frac{T_{\text{round}}}{2} \right)$$

$$\text{Skew} \leq \left(\frac{T_{\text{round}}}{2} \right) - \text{min}$$

$$\leq \left(\frac{T_{\text{round}}}{2} \right) \quad \begin{matrix} \text{minimum 1 way network} \\ \text{delay (\approx 0)} \end{matrix}$$

(ii) Berkeley Algo = only internal synchronisation

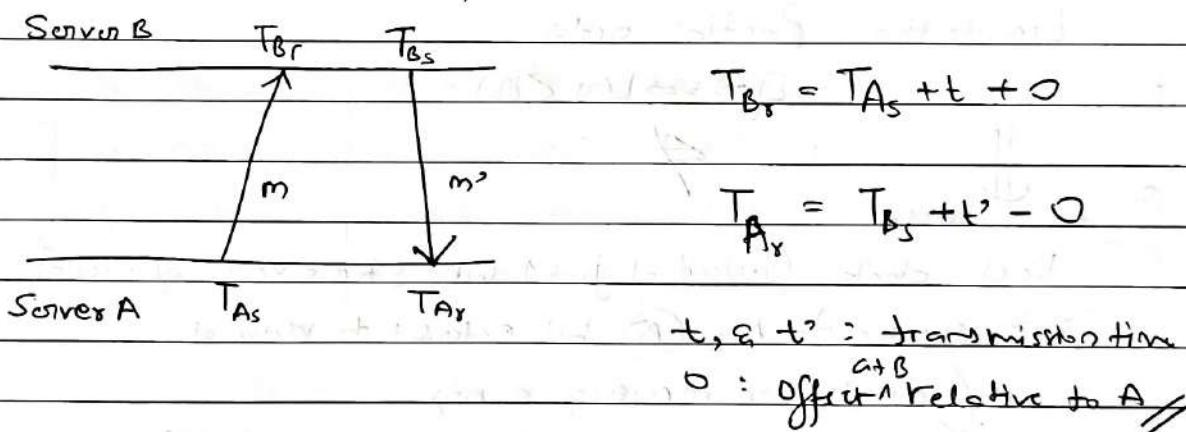
- Server asks all their time
- Christian algo on each
- Avg the local time



Send the offset to all

(iii) Network Time Protocol

- Multicast timestamps - low accuracy
 - procedure call (christian algo) - high accuracy
 - Symmetric mode = high accuracy
- ?
- \Rightarrow • A & B exchange messages
- Record the send & receive timestamps
 - calculate the offset b/w them



Normal clock req synchronization

L4



logical clocks : time sync not req

order of events req

- Happened Before Relationship

(HB) : denoted by \rightarrow $e \rightarrow e'$ means e happened before e' $e \rightarrow_i e'$ means e happened before e' as observed by p_i \rightarrow if $\exists p_i, e \rightarrow_i e'$ then $e \rightarrow e' \quad \{ \Rightarrow l(e) < l(e') \}$

(X) true

as $e \rightarrow e'$ or $e' \rightarrow e$

Q

NOTE: Algorithm (Lamport's Timestamping)

(A)

(i) $\text{ft} L_i = 0$

(ii) increment on any event

(iii) on send message \rightarrow send L_i (iv) receiver: $\text{ft} L_i = \max(L_i, t)$

~~~~~ received

(v) then increment  $L_i$  before timestamping receive msg

Limitation: Partial order

 $\text{ft} L_a < \text{ft} L_b$ ↓  
hence,

Vector clock (instead of just 1 clock, store a vector of clocks)

⇒ similar to Alg (A) but extended to vector of

(clock) on receiving a msg

note: increment just own (but take max for all)

$$VT_1 = VT_2 \Rightarrow VT_1[i] = VT_2[i] \quad \forall i$$

$$VT_1 \leq VT_2 \Rightarrow VT_1[i] \leq VT_2[i] \quad \forall i$$

$$VT_1 < VT_2 \Rightarrow VT_1[i] < VT_2[i] \quad \forall i$$

$$\text{eg } VT_1[j] < VT_2[j] \quad \exists j$$

$\Rightarrow VT_1$  is concurrent with  $VT_2$

iff  $(\text{not}(VT_1 \leq VT_2) \wedge \text{not}(VT_2 \leq VT_1))$

{ NOTE: ordering of events based on timestamp  
is  $VT_1 < VT_2$  } \*

Causality Violation: occurs when order of messages cause an action based on information that another host has not yet received

↓

$\Rightarrow$  If incoming vector timestamp  $<$  local V.T.S

# Global State ( $\equiv$  state of all processes + state of all communication channels)

- includes instantaneous state of each process

how to get it without synchronising the clocks?

- 3 types of events: (i) local computation

(ii) sending a message

(iii) receiving a message

$e_i^n$ : nth event at  $p_i$

NOTE:

(i) History ( $p_i$ )  $\Rightarrow h_i = \langle e_i^0, e_i^1, \dots \rangle$

(ii) prefix history ( $p_i^k$ )  $= h_i^k = \langle e_i^0, e_i^1, \dots, e_i^k \rangle$

$S_i^k$ :  $p_i$ 's state immediately after kth event

(iii) For a set of processes  $\langle p_1, p_2, \dots, p_n \rangle$

$$H = \bigcup_i (h_i) \quad \{ \text{global history} \}$$

NOTE: a cut  $C \subseteq H = h_1^{c_1} \cup h_2^{c_2} \cup h_3^{c_3} \dots h_n^{c_n}$

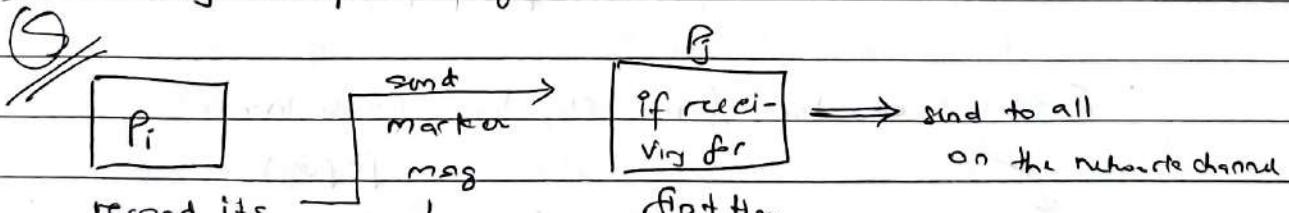
Also frontier of  $C = \{e_i^{c_i}, i=1, \dots, n\}$

global state  $s$  corresponds to cut  $C = \bigcup_i (S_i^{c_i})$

consistent: iff  $\forall e \in C ( \text{if } f \rightarrow e \text{ then } f \in C)$

consistent iff it corresponds to a consistent cut!

(L-5) Chandy Lamport Algo



Own state

(Start recording all)

msg in that channel until received a marker

(incoming)

• ensure consistent cut? Proof by contradiction

• channel state?

(delivery of marker E a message): FIFO

# Run - a total ordering of events in  $H$  that is consistent with each  $b_i$ 's ordering

Linearisation - run consistent with happens-before ( $\rightarrow$ ) relation in  $H$

It passes through consistent global states

NOTE: A global state  $S_k$  reachable from  $S_i$  if there is a linearisation that passes through  $S_i$  & then through  $S_k$

The d.s. evolves as a series of transitions between global states  $S_0, S_1, \dots$

### (L-6)

- Liveness: something good will happen eventually

$\text{Liveness}(P(S_0)) = \forall L \in \text{linearisations from } S_0, L \text{ passes through } S_L \& P(S_L) = \text{true}$

State  $S_0$  satisfies liveness property

- Safety: something bad will never happen

$\text{Safety}(P(S_0)) = \forall S \text{ reachable from } S_0, P(S) = \text{true}$

State  $S_0$  satisfies the safety property

- Stability: once true, forever stays true (stable liveness property)

(if a state once false, always false (stable safety property))

is true, all

the states following it will remain true)!

Distributed Debugging: to a monitor, all global states sent

# Possibly  $\phi$ : There is a consistent global state  $S$  through which a linearisation of  $H$  passes such that  $\phi(S)$  is True

Definitely  $\phi$ : for all linearisations  $L$  of  $H$ , there is a consistent global state  $S$  through which  $L$  passes such that  $\phi(S)$  is True

NOTE:  $S$  is a consistent global state iff

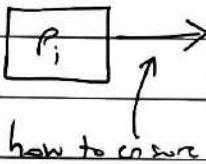
$$\forall [s_i[ij] \geq \forall [s_j[ij]] \quad \forall i, j = 1, 2, \dots, N$$

NOTE:  $S = [s_1, s_2, \dots, s_n]$

This means no. of  $p_i$ 's events known at  $p_j$  when it sees  $s_j$   
Should not be more than no. of events that have occurred at  
 $p_i$  when it sees  $s_j$

(L-#)

Check condition again!!



Reliability?

Orderedness?

M1 Reliability } II) B-Multicast

Only for multicast it should have happened before relationship

How to achieve ordering

(i) FIFO

\*\* ↑↑

-(ii) causal

FIFO ← causal

(iii) Total

M2 Reliability } II) R-Multicast

Integrity: a correct process p delivers a message at most once

Trusted ↗ MSR

(multicast)

(no retrans)

# overall liveness

What if sender fails → reconstruct

↓

Overhead

↓

B-Gossip

tend to B

random

targets)

# implementation

→ sequence vector

FIFO } Multicast

Causal

eventually delivery

TOTAL Order } → delivery of message

↳ 2 ways: central server / 1SIS

↳ proof

(N) X (K)

K

R

B

$$SK = \boxed{M}$$

(M) X (N)

□ -&gt; M

MXN

□ -&gt; M

□ -&gt; M

L-8

Mutual Exclusion - in single OS : Semaphore, Mutex, Conditional

Safety



Variables, Monitor

liveness



in DS - central superv algo

Ordering



Ring based "

↳ as in how it was sent (not received)

ricart-agrawala " - Causality &amp; Multicast

mackawa "

" - no liveness

## Performance eval metrics

(i) Bandwidth - total no. of messages sent in each enter &amp;

exit operation

incurred

(ii) Client delay - delay <sup>↑</sup> by a process at each enter &

exist operation (When no other process is in CS or waiting)

(iii) Synchronization delay - time interval bet' one process

exiting the CS &amp; the next process entering it (when there is only one process waiting)

measures the throughput of the system

L-9

## Leader Election

→ algorithm should guarantee: Safety

liveness

↳ Ring election algo

Bully algo → only to seniors (if no replies declare itself)

↳ Performance Analysis

do wait still if no → rerun!!

(i) Bandwidth usage: Total no. of messages sent

(ii) Turnaround time: The no. of serialised message transmission time between the initiation &amp; termination of a single run of the algorithm

# consensus impossible in synchronous system  
(so is election!)

$\rightarrow$  Reliable & Totally Ordered Multicast

(L-10) Consensus  $\equiv$  Interactive Consistency Problem  $\equiv$  BGZ

E11  $\nwarrow$  for synchronous: soln exist

- Properties: (i) Termination

- (ii) Agreement

- (iii) Integrity

(for synchronous system: fail? exist)

Synchronous: Send to all & wait for  $\infty + (f+1) * T$  for upto f failed

(Process  
(after  $f+1$ ) rounds, algorithm will stop)

$$\begin{array}{c} f \geq [N-1] \\ \diagup \quad \diagdown \\ 2 \end{array} \quad \boxed{3f < N}$$

(Proof by contradiction)

Byzantine Problem

$\rightarrow$  consensus safety

Asynchronous: Paxos }  $\rightarrow$  either safety or liveness  
Raft }  $\rightarrow$  (used for deciding a single value)

(L-11) Paxos  $\rightarrow$  proposer + acceptor + learner

What is Paxos - for single value  
if multiple  $\Rightarrow$  LogConsensus (Raft)

(L-12) Raft

Objective: Whenever possible, select the alternative that is the easiest to understand

leader based!

Candidate

followers

(L-13)

- RPC - Invocation semantics - atleast once, atmost once, many times!
- Marshalling & unmarshalling (as each system have its own format)

$\downarrow$  to EDR       $\uparrow$  from EDR

$\rightarrow$  external data representation

(to/from) (platform dependent)

- all or nothing principle  
 → seq of op - by Client - transform Server data from  
 (L-16 § 17) → every step is an RPC one consistent state to  
 Transactions → ACID properties another  
 ↗ isolation major issue

option 1: serial execution {concurrency reduced}

option 2: serial equivalence : criterion for correct

i.e.  $T_1 = \{op_1, op_2, op_3\}$  concurrent execution

$$T_2 = \{op_4, op_5, op_6\}$$

$$O = \{op_1, op_2, op_4, op_3, op_5, op_6\}$$

is serially equivalent if,

result same as either

$$\{op_1, op_2, op_3, op_4, op_5, op_6\}$$

or

$$\{op_4, op_5, op_6, op_1, op_2, op_3\}$$

• conflicting operations : read & write

↑  
write & read

} on same object

in serial equivalent,

write & write

all the conflicting

read & read (X)

operations will have

an order → i.e. either ( $T_1 \sqsubseteq T_2$ )

or ( $T_2 \sqsubseteq T_1$ )

# ways to achieve concurrency control

pragmatic → a) Locking

optimistic → b) Optimistic CC

c) Timestamp ordering \*

must commit / abort!

• Read → ↗ write

must commit / abort!

• Write → ↗ Read / write

→ Not serial equivalent  $\Rightarrow$  2 Phase Locking

NOTE: in read-walk lock, head  $\rightarrow$  walk promotion

happens only if both operations in same transaction!!  
(otherwise blocked!!)

{gives signal equivalence & locks} (Proof)

# 2 phase locking  $\rightarrow$  Strict 2 phase includes release locks only

(i) growing phase  $\curvearrowright$  at commit point)

(ii) Shrinking phase  $\rightarrow$  acquiring  $E_1$  promoter

(ii) Shrinking phase  $\rightarrow$  acquiring  $E_1$  promotion

downside: locking !!  $\rightarrow$  deadlock

<sup>measuring</sup>  
3<sup>rd</sup> cond<sup>n</sup>S: (i) exclusive lock modes

$b = \sum t_i b_{i+1} d_i$ ?

(ii) ~~lecter~~ no preemption

- Local all objects at start (iii) Circular wait in wait for graph (iv) Hold & wait tool

(iv) Hold & wait) too!

- lacks timeout
- deadlock detection: keep track of wait for graph & find cycles  
Init  $\Rightarrow$  breaks the (iii) point

(L-18)

## Distributed Transaction reg Coordinator -

↳ How to commit?

→ an issue :- One phase commit } - not so effective (how many?)

• 2 phase commit - coordination

(phantom deadlocks) ← { how to avoid deadlocks? ← → as locks not released

## How Recovery?

- Recovery manager
  - Recovery file

on detection,  
lowest priority transaction aborted to break the cycle!

until commit or abort

→ Centralised detection: global graph

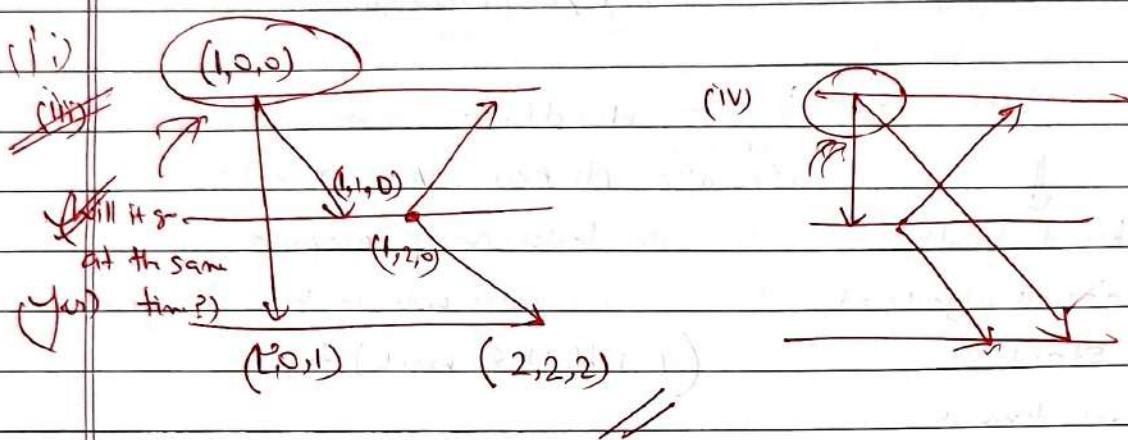
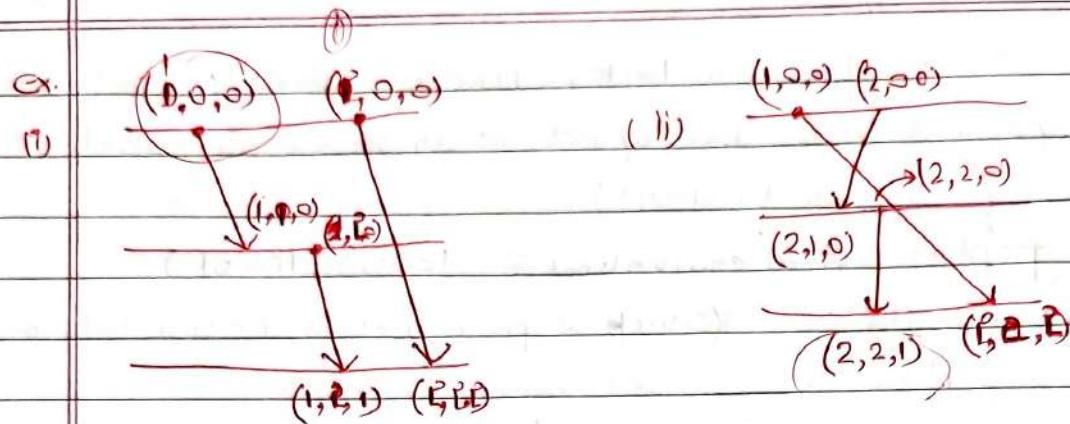
→ decentralised deadlock detection

↳ edge chasing

(if a server notices a transaction

Waiting for another transaction and

It waiting for another =  
by sending protreq) 1



(Lec 16 - 17)

## ° Optimistic Concurrency Control

a) 3-phases

working  $\rightarrow$  validation  $\rightarrow$  update phase $\downarrow$ 

Stores

tentative

version

(each obj)

Validated

to see if

any conflicts

with other

transaction

if no conflict, result/

transactions made

is updated

(otherwise)

## b) validation rules

|           |               |   |
|-----------|---------------|---|
| until     | read - read   | } |
| committed | write - read  |   |
|           | write - write |   |

$\rightarrow$  done by backward or forward validations

### c) Backward validation

Read set of  $T_V$  compared with write set of  $T_i$

prev transaction

$i \leftarrow V$

### d) Forward validation

Write set of  $T_V$  compared with read set of  $T_i$

$i \leftarrow V$

NOTE: Write-write conflicts don't matter for OCC validation

because they don't cause stale read or break serial order, it's just that the last writer wins.

### • Timestamped Ordering

a) each transaction assigned an id

b) if write  $\rightarrow$  check <sup>last</sup> read & write on that

obj is less than incoming one  
if read  $\rightarrow$  check last <sup>write</sup> read is lower id

c) never results in deadlocks

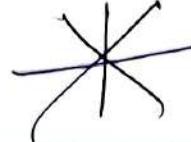
### d) ordering rules

write read - write if  $\geq$  read max timestamp

write write ?

read write

} similar conditions



d) T.O : Write rule

if  $T_c \geq \max$  read T.S on D

{ }  $T_c \geq \text{write timestamp on committed version of D}$

Q  $\rightarrow$  list of tentative writes sorted

if  $T_c$  in T.W list  $\rightarrow$  update by ids(asc)  
else add it

else

// drop

e) T.O : Read rule

if  $T_c > \text{write T.S on committed version of D}$

$D_s = \max \text{ write T.S} < T_c$

if ( $D_s$  committed)  $\rightarrow$  add  $T_c$  to RES list

read  $D_s$

list of ids that

have read committed value

else if  $D_s$  written by  $T_c$   $\rightarrow$  read  $D_s$

else wait till  $D_s$  committed or aborted then  
reapply rule

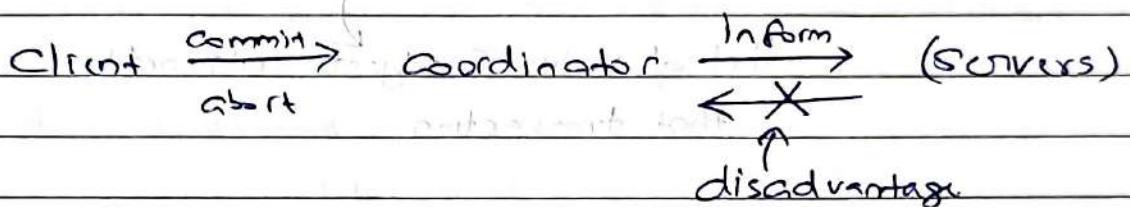
3 else {

// abort

(U8-20)

either all or nothing

- o Distributed Transaction Atomicity == Consensus
- o requires coordinator : to whom req is fwd  
other servers serve as participants
- o It's an asynchronous system !!  
→Req safety property = never have some agreeing to commit and others agreeing to abort
- o First cut: One phase commit



→ blocked until a failed coordinator recovers

- o Two phase commit Protocol → participants save state in mem.
- Phase 1: Collect all commits → tells to Coordinator logs
- Phase 2: If all replies commit → tells to Coordinator logs  
else sends abort message

# how to handle with loss

(i) participant crashes \*

(ii) can't commit? loss ? Participant after a timeout eventually says 'No'

(iii) yes or no loss ? aborts transactions Coordinator after some timeout

(iv) do commit loss: after saying yes, it cannot abort by itself → has to send ~~get Decision~~ <sup>get Decision</sup> → retry !!.

- Phantom deadlock

→ false detection of deadlocks that don't actually exists !! ))

as edge chasing messages may have stale data

→ leads to spurious aborts

- Recovery

→ maintains recovery file consisting of object, transaction status, Intentions list

list of values & objects altered by that transaction

- Logging

→ contains history of all transactions performed by a server

→ this is filled during commit or abort

prepares to commit

first transaction

at that time, objects added in Intentions list  
then commits → status updated

Log → Checkpoint : saved snapshot of the database

Updates or system at a certain moment

Preparation

Commit

Why req?

else need to replay all

log from the start

→ by reading log

- After a crash, 2 ways to recover:

(i) Forward Recovery (Redo based)

→ Start at the last checkpoint & move forward

In the log

[that is committed]

→ (like literally apply all the log events)

(ii) Backward Recovery (Undo based)

→ Start from the end of the log & move backwards

→ i.e. undo committed transactions' updates

- Shadow Versions

→ make a copy of the data when committing

instead of modifying the original

→ during prepare phase - creates a shadow version

& on commit phase it replaces the old map

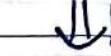
with new one

(Logs 24 - 26) → 23 22 21 20 19 18 17

### a) Security Threats

- Leakage - Confidentiality needs to be ensured
- Tampering - Integrity
- Vandalism - DOS - Authority
- and many other factors like Availability

b) Security Policy indicates which actions each entity (user, data, service) is allowed or prohibited to take



Security mechanism implements and enforces the policy

### c) Cryptography

- Private/ Shared Key Cryptography

→ how to Share Key?

→ key length: 128 bits

- Public Key Cryptography

→ key length: 512, 1024 bits etc

### d) Encryption

→ Intruder Types

- Passive - only listens to encrypted texts
- Active - can alter messages

→ Substitution Cipher

DES: Data Encryption Standard

◦ encrypts data in 64 bit chunks

◦ Key length: 56 bits (actually 64 bits)

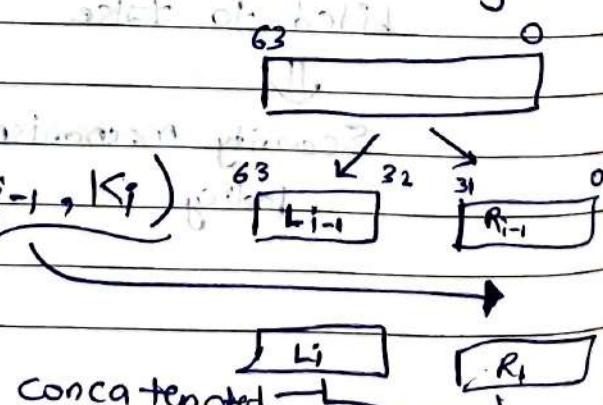
◦ It's a 16-round Feistel cipher - 19 stages of transformations (for parity)

~~Each of~~ each round has its own keys

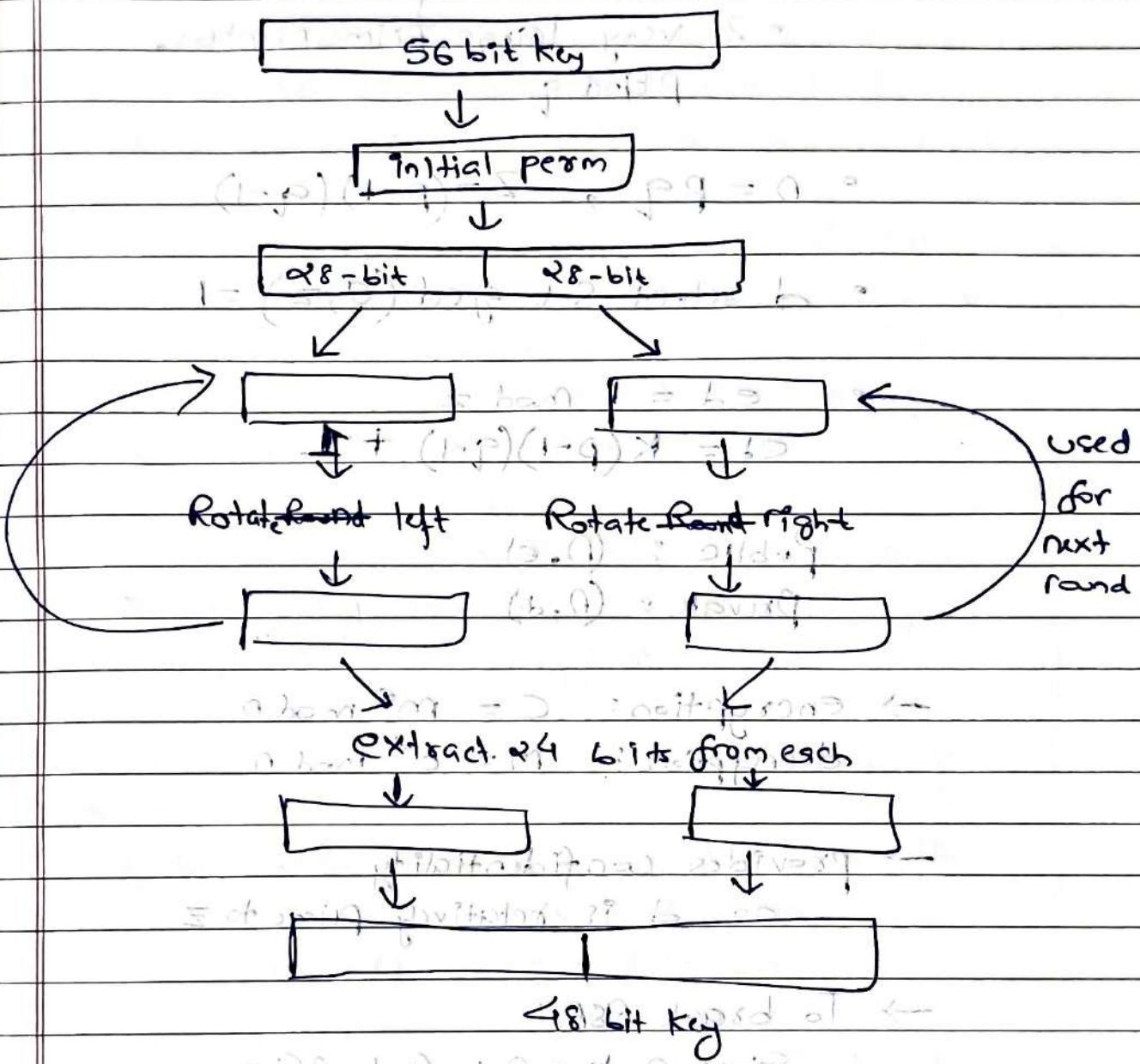
~~Each of~~ generated from original key: (48-bit key used)

$$L_i = R_{i-1}$$

$$R_i = L_{i-1} \oplus F(R_{i-1}, K_i)$$



° Key generation using 56 bit key



° provides confidentiality - (56! : Reason)

Under this only

→ e) Key Distribution Problem

f) Trusted centralized Agent

→ only need to know about trusted agents' key

↳ ~~↳~~ ~~↳~~ ~~↳~~ ~~↳~~ ~~↳~~ ~~↳~~

↳ ~~↳~~ ~~↳~~ ~~↳~~ ~~↳~~ ~~↳~~ ~~↳~~

↳ ~~↳~~ ~~↳~~ ~~↳~~ ~~↳~~ ~~↳~~ ~~↳~~

## Public Key Cryptosystems : RSA

→ how to generate key

- 2 very large prime numbers,  $p \neq q$

$$\bullet n = pq, z = (p-1)(q-1)$$

- $d$  select s.t  $\gcd(d, z) = 1$

$$\bullet ed = 1 \pmod{z}$$

$$ed = k(p-1)(q-1) + 1$$

Public :  $(n, e)$

Private :  $(n, d)$

→ encryption:  $c = m^e \pmod{n}$

decryption:  $m = c^d \pmod{n}$

→ provides confidentiality

as  $d$  is relatively prime to  $z$

→ To break RSA:

Since  $n$  known: find  $p \& q$

### Hash functions

(i) Secure Digest functions

$$h = H(M)$$

provided,  $h \neq M$

a) Given  $M$ , it is easy to compute  $h$

b) Given  $h$ , it is hard to compute  $M$

c) Given  $M$ , it is hard to find another message  $M'$ , such that  $H(M) = H(M')$

## example MD5

Input broken into 448 bit blocks

padded to turn it into 512 bit blocks

Message digest size: 128 bits

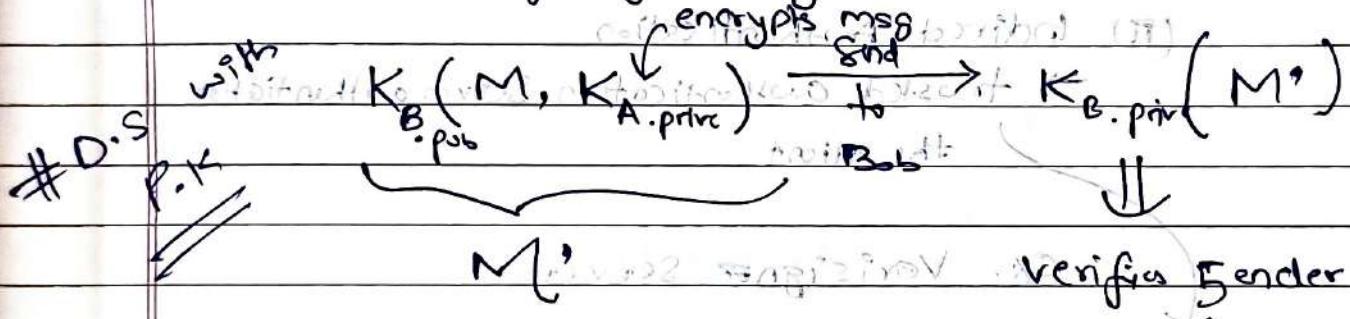
Protocol ← done through PKC

## ii) Digital Signature

- to protect sender receiver from each other
- Sender encrypts the message using their own private key

And receiver can verify it using public key

NOTE: If confidentiality needed



→ provides data integrity  
message authentication  
non-repudiation

→ generally if Unsecure channel

→ PKC

else if secure Channel

→ establish the shared secret key

Using hybrid method

→ use it to produce low cost signatures

(a short piece of

info generated using  
symmetric keys)

- Message Authentication Code.

→ provides integrity & message auth

but not non-repudiation (!!!)

as signature ← as if we can verify a MAC, we can

Share(?)

also create it

→ proves identity using a trusted CA

### iii) Digital Certificates

→ a document containing a statement

(usually short) signed by a principal

→ binds a public key to a real identity

for Pt's proof

#### o Authentication

##### (i) Direct Authentication

→ Server and clients contact

→ behind firewall

##### (ii) Indirect Authentication

→ trusted authentication. Server authenticates the client

#### o CX. Verisign Services

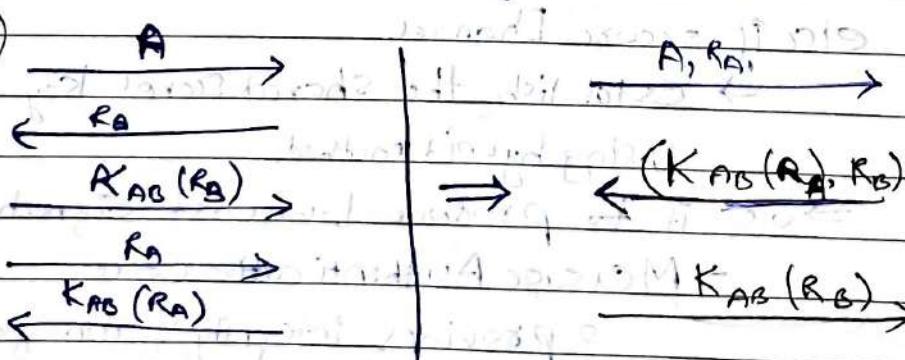
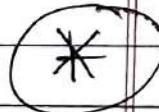
# Solution based on Secret Key Cryptography

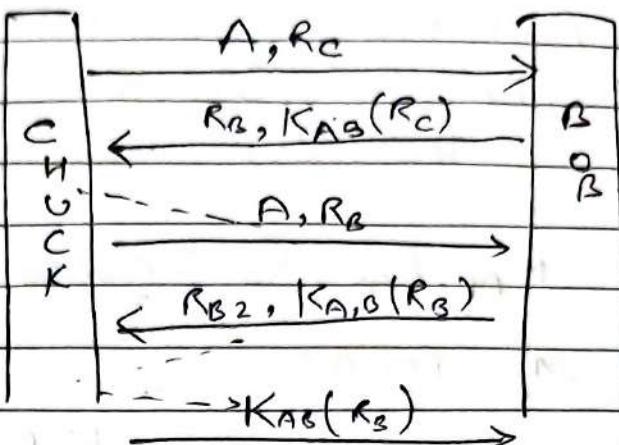
\* → 3-way handshaking

→ Trusted third party (KDC)

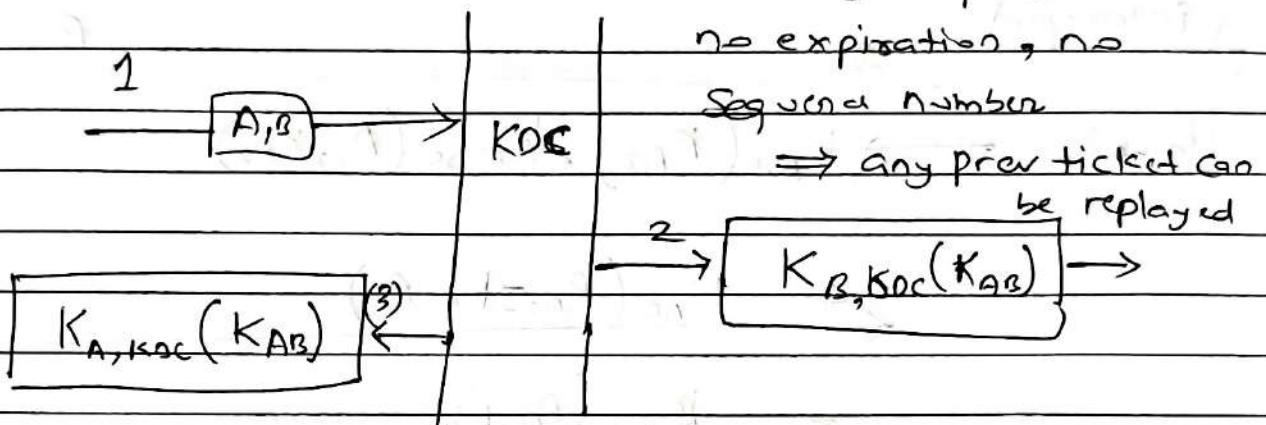
# Soln based on Public Key Cryptography

→ Public Key Authentication



ReplayReflection  
Attack

based

# KDC<sup>n</sup> Protocol  $\Rightarrow$  no timestamps, no nonce,  
no challenge response.

Pb: A may communicate to B even before B got the key from KDC  
so in (1)

~~$K_{A,KDC}(K_{AB}), K_{B,KDC}(K_{AB})$~~  ticket

$\Rightarrow$  let Alice maintain setting up a

connection to Bob

(vulnerable) if  $K_{B,KDC}(K_{AB})$  old released

Send to B:  $(A, K_{B,KDC}(K_{AB}))$

It can establish

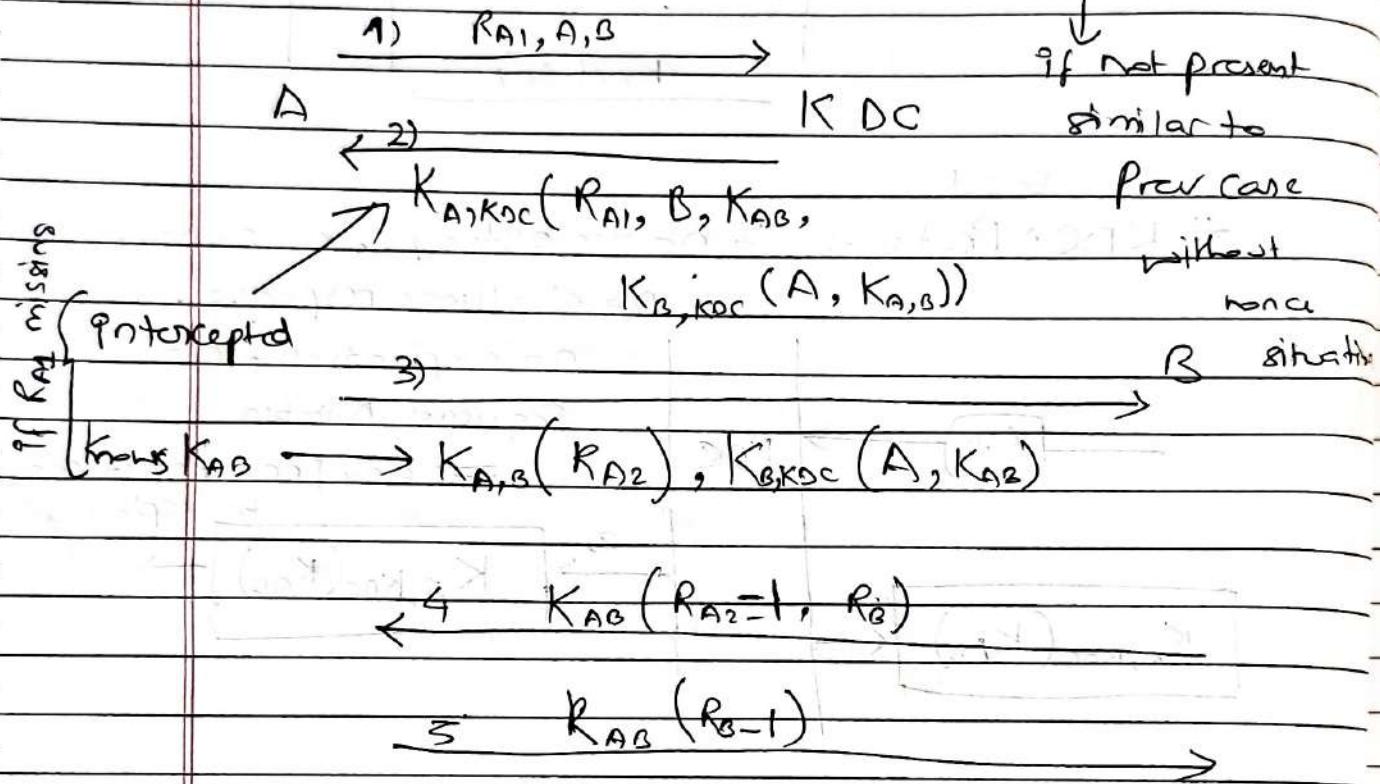
a session as no verification!

for session establishment

## Authentication Using KDC

(Needham Schröeder Protocol)

$R_{A1}, R_{A2}, R_B$  : nonces = purpose of freshness



Vulnerable to replay if  $K_{AB}$  (released)

- Case 1 :
- If he knows  $K_{B,KDC}$  old without  $R_{A1}$  present
  - Then he can get  $K_{AB}$ .

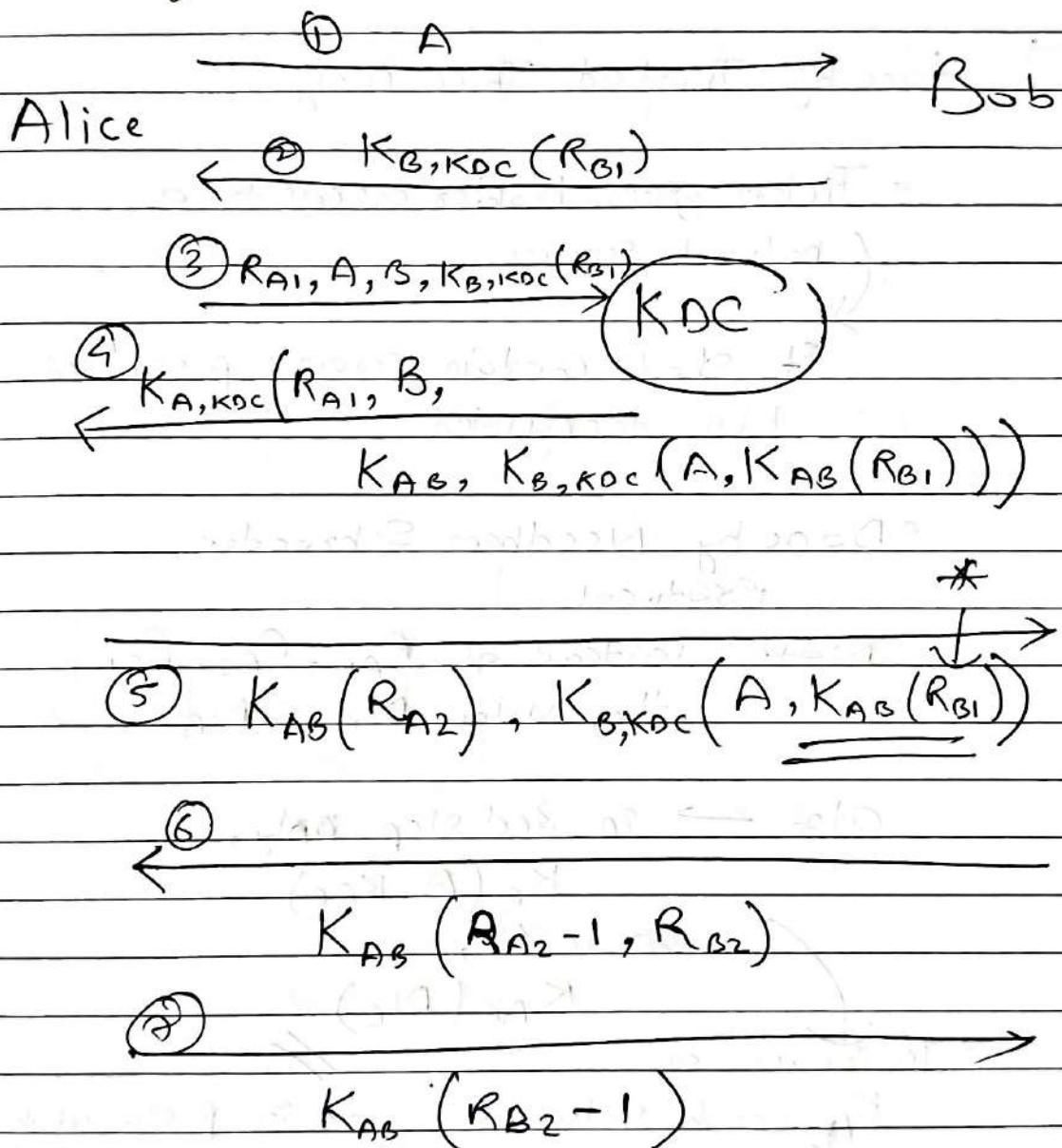
Case 2 : If  $B$  missing → replay attacks can happen

Identifies who  
he wants to talk

Case 3:  $K_{AB}^{\text{old}}$  is released

- ↳ old message can be replayed
- and session establishment

Defense



## # Authenticating to Multiple Servers

- Ticket based system

(Else Vulnerable)

↳ as in if every server knows  
the password  $\Rightarrow$  one broken, all unsafe

done by Trusted Third Party

- Ticket gives holder access to a

(network service)

It should contain server's password  
but encrypted

- Done by Needham Schröeder  
Protocol

Note: instead of  $R_A1, R_A2, R_B1$   
they using  $N_A \oplus N_B$

also  $\rightarrow$  in 3rd step only,

$K_B(A, K_{AB})$

so in 4th,

$K_{AB}(N_B)$

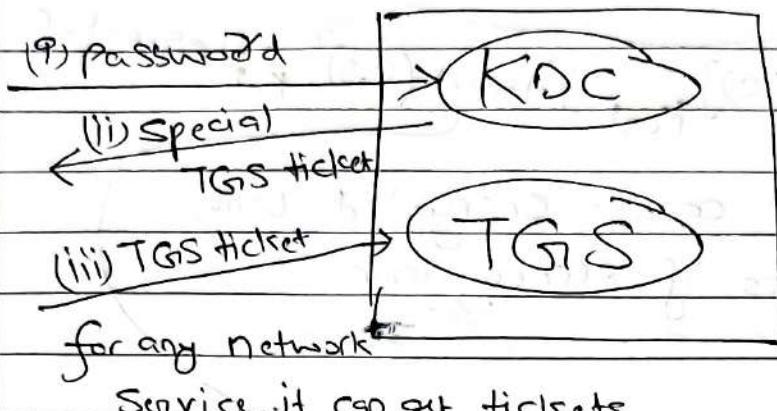
- Insecure as

$K_{AB}$  can be stolen

as in password to  
authentication/ticket

- also need to send password each time generate  
for ticket generation,

## # Better Solution Two Step Authentication

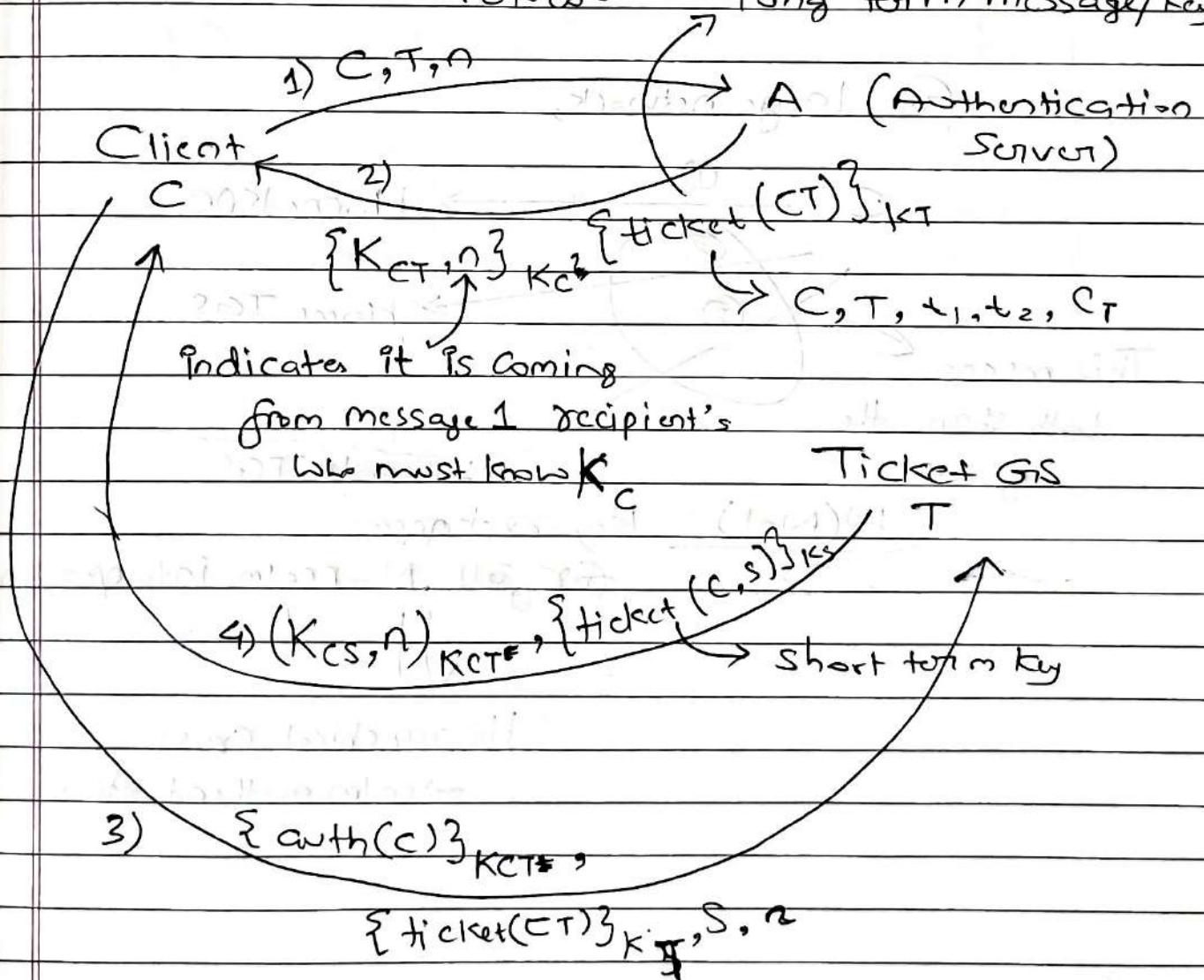


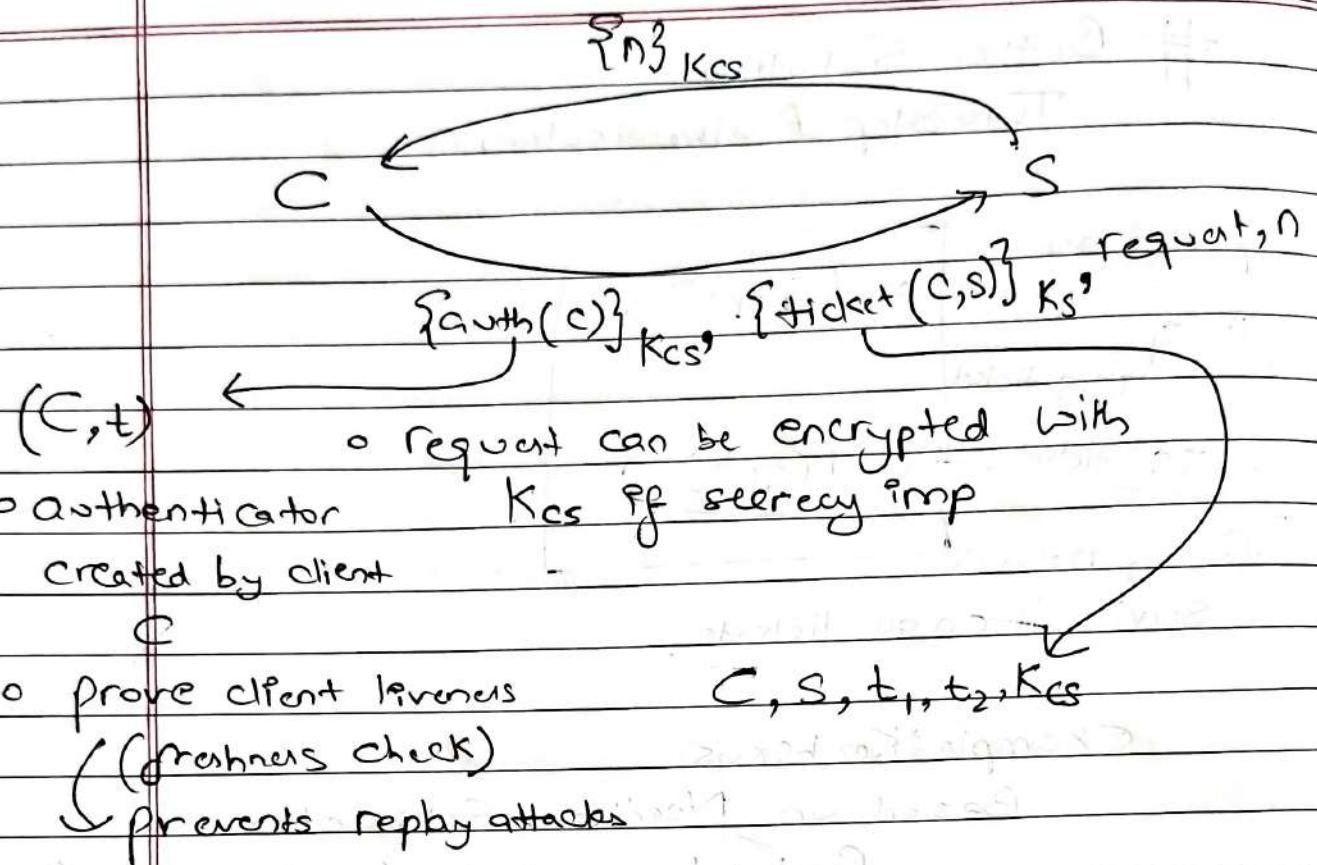
Example: Kerberos

Based on Needham-Schroeder

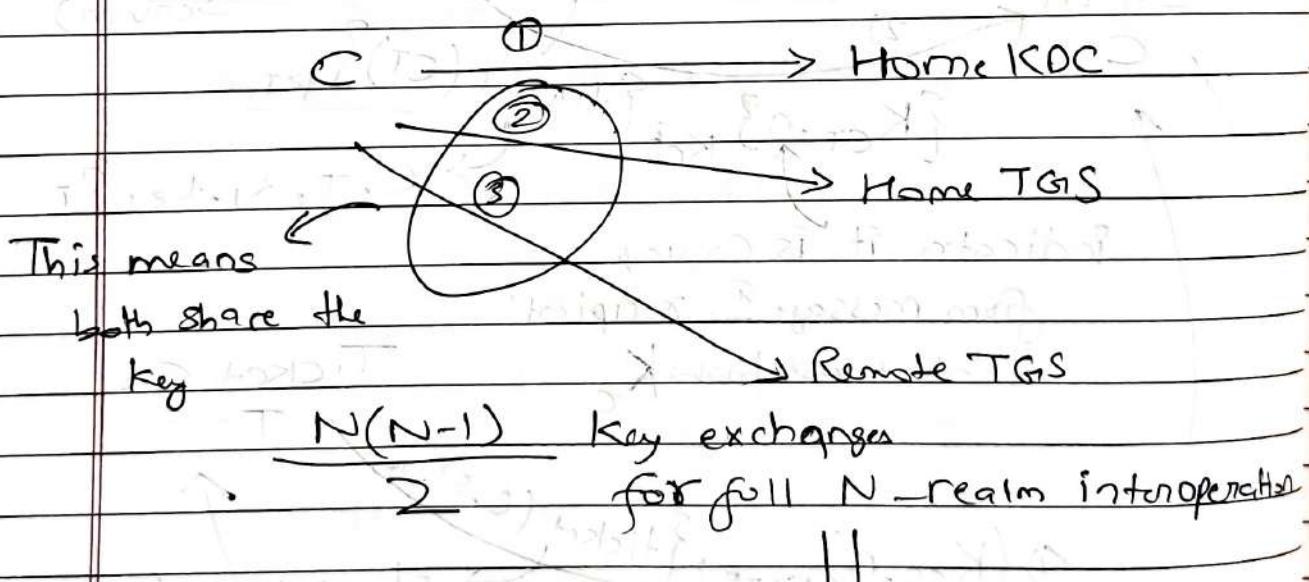
Protocol

long term message/key





for large network,



Hierarchical cross

-realm authentication

(Video)!

(users)

→ not just tcp, SMTP over SSL

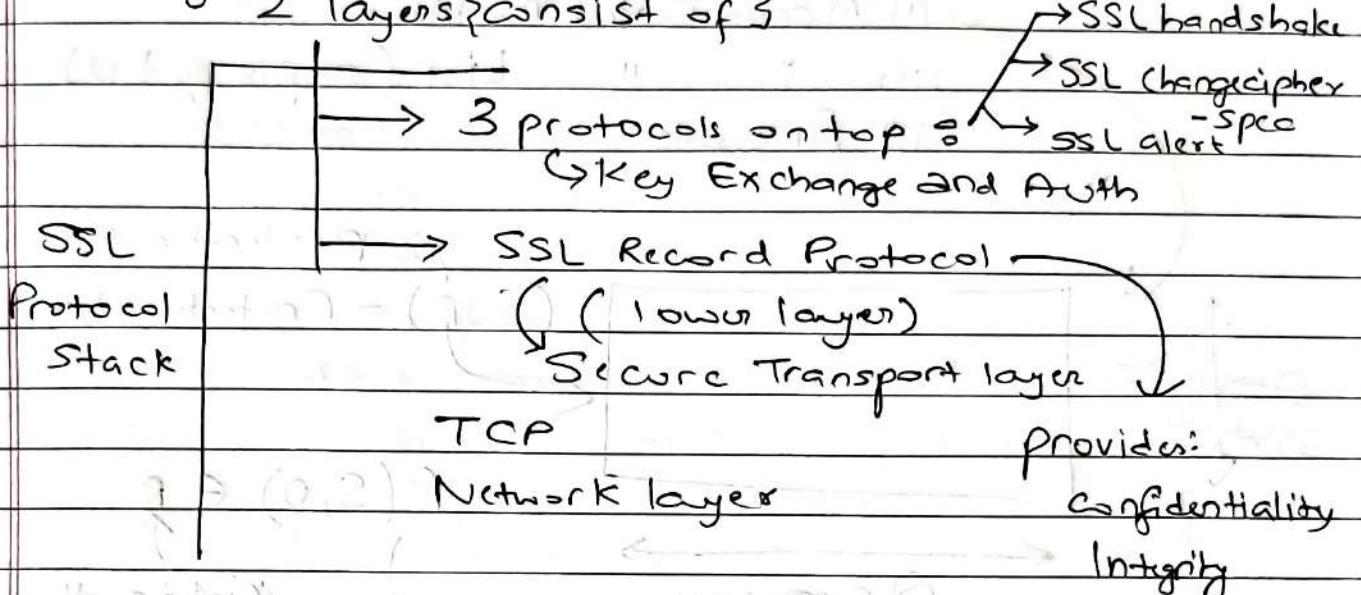
## # Secure Socket Layers

(465 port no)

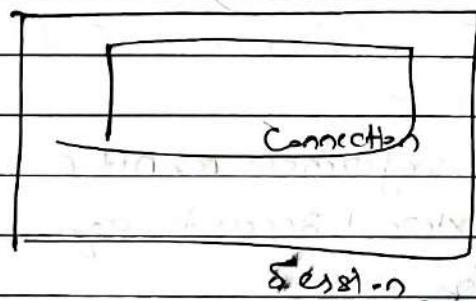
layer on top of tcp

- tcp + SSL - authenticated, reliable, end-to-end, encrypted stream

- 2 layers consist of 3



- provides confidentiality, integrity and authentication



(# Read one

more time

(for SSL)

Specific to each layer

L24-26: (Slide 86 - 108)

→ user req, smtg  
 ↓

to verify if can access  $\Rightarrow$  do

access  
control

## # Access Control <

→ verifying access rights

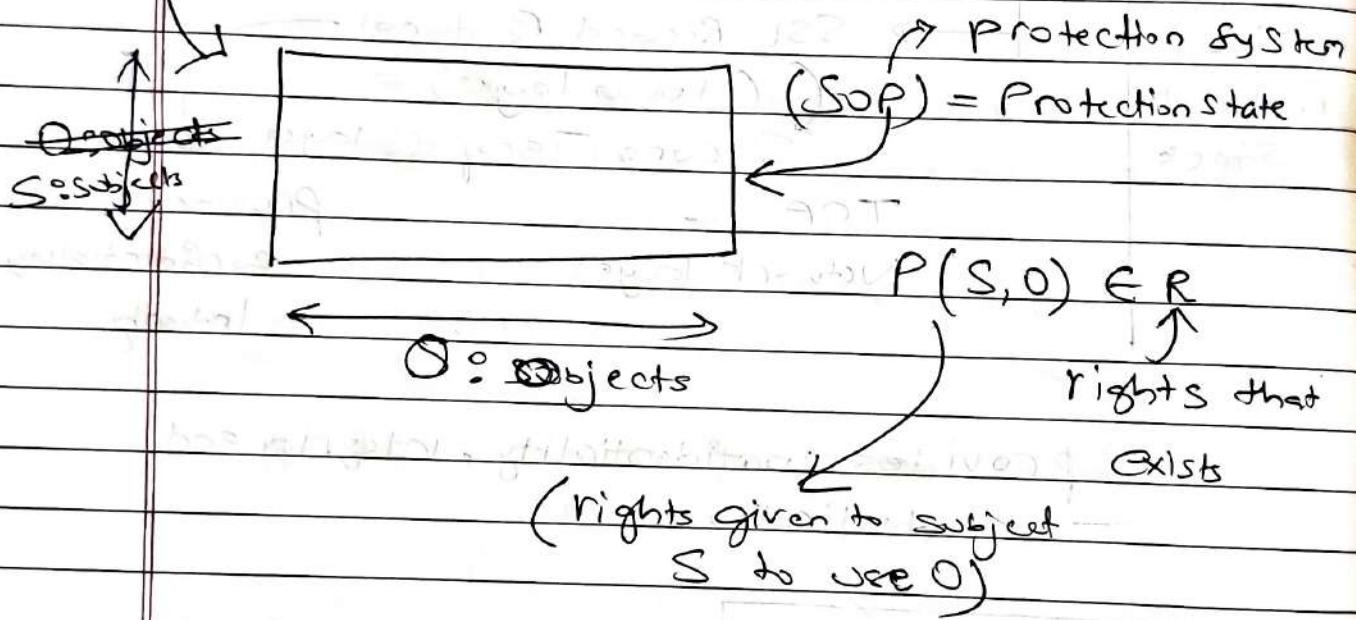
→ granting " " " = authorization

→ access control models

(i) Access Control Matrix

(ii) " " List (capability Lists)

(iii) Firewalls



## Firewalls

→ Special kind of reference monitor  
to control external access to any part of a d.s

→ Models: (i) Packet filtering gateway  
(ii) Proxy

(Lec 27 - 28)

### a) Key-value Abstraction

→ a dictionary data structure

→ examples: SQL

## NoSQL

## b) Distributed Hash Table (DHT)

→ a distributed key value store

→ Uses a hash function to map

Key → node

hash(file)    value(file)

→ different protocols to implement DHTs

example. Chord, CAN, Pastry,

## Tapestry

### c) CHORD

→ properties: Load balance

## Decentralized

Scalable - Cost (lookup) =

$$C(\log N)$$

## High availability

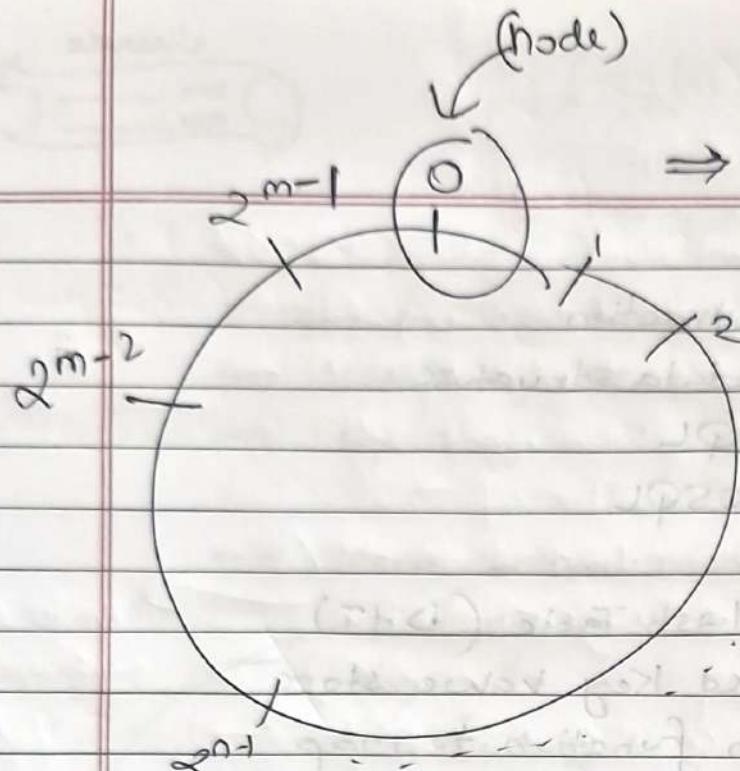
## Flexible naming

$(O(n \cdot m))$  } Every node gets an ID? using hashmap  
" key " " " "

→ All IDs lie on circle (ring)

→ a key is stored on the first node clockwise from the key's ID → such that it's  $\geq k$ ,

i.e called successor of  
key K



classmate  
Date \_\_\_\_\_  
Page \_\_\_\_\_

⇒ NOTE: m bit is the d/p  
of the address space  
in ring

- find the id which  $\geq k$   
i.e. successor of  
key  $k$   
( $\therefore k \% 2^m$ )
- if none present,  
go to smallest ID  
available

⇒ This is Consistent hashing

⇒ Useful: When a node enters or leaves

contents need to move

exists only  $\frac{1}{n}(k)$  of them shifts

( $n = \text{no. of nodes}$ )

NOTE: key mod  $2^m$  stores the data we want to store

Note:  $2^m$  possible locations for storing a node but only certain locations actually have a server (called node) !!

- That's why we search for Successor(Key) to store the value

→ helps in load balancing

→ Consistent Hashing

With  $K$  keys and  $N$  peers, each peer stores

$O(\frac{K}{N})$  keys (i.e.  $\leq C \cdot \frac{K}{N}$  for some constant  $C$ )

# How to perform lookups

$$(i) O(n)(T-C) + O(1) S \cdot C$$

Succesor  $\leftarrow$  initial value (current node / queried node)

(A)

i. Locate successor ( $key$ ) where  $key \neq i$ :  
 if ( $key \in (i, \text{successor}]$  then  
 return (successor))

\* else return successor. Locate - successor( $key$ )

(ii) Finger Table

- each node aware of  $\approx m$  other nodes
- maintains a finger table with  $m$  entries

i<sup>th</sup> entry of node  $n$ 's finger table =

$$\text{successor}(n + 2^i)$$

$i$  ranges from 0 to  $m-1$

Algorithm for look up!

→ (A) remains same except (\*) part gets updated:

else forward to the largest finger entry  $\leq K$

NOTE: Every query is for key  $K$ !!

$$\text{Total cost} = O(\log n)$$

optimize routing

NOTE: finger table entries used for both  
lookup, insertion, deletion etc  
as long as  
the finger table & successor correct!

Also each node stores a succesor list

Pb of all  
node dying  
 $(\frac{1}{2})^r$

$\hookrightarrow$  size  $\propto (= 2 \log N)$   
(nearest & preserve  
successors) ring successor

$\Rightarrow$  at least 1

alive:

$1 - (\frac{1}{2})^r$

# Failure Recovery

NOTE: Successor list helps here by replacing  
it with first live entry in the list

$\Rightarrow$  single alive  
half alive

$(1 - (\frac{1}{2})^r)^{\frac{N}{2}}$

Stabilize will correct finger tables &nd  
successor list entries pointing to failed

node

Storing keys

NOTE: just not the successors stored  
we replicate key value pair at & successor  
& predecessors!

# New nodes joining!

$\rightarrow$  simple: Add the node

& then update finger tables

$\rightarrow$  no. of finger tables affected due to new  
node addition =  $\log N$

→ No. of messages per node join (to initialise the new node's finger table)

$$= O(\underline{\log(n)} * \underline{\log(n)})$$

↑  
no. of entries

↑  
each req

lookup

## # Concurrent Joins

i) Updating every finger table perfectly is expensive

ii) Chord lookup — doesn't require every finger table to be correct

iii) Reg. only 2 invariants:  $\rightarrow \text{Successor}(n)$

Invariant 1: each node must have the successor pointer — i.e. first entry in the finger table

Invariant 2: The correct successor of each key  $K$  must exist somewhere on the ring

$\rightarrow \text{Successor}(1 <)$

## # Stabilization Protocol

a) When a node joins

$$\text{next}(n) = \text{successor}(n)$$

$$\text{notify}(\text{next}(n))$$

} joining node notifies its successor

b) When a node gets notified  $\therefore n \leftarrow \text{notify}(n')$   
 node updates If  $\text{prev}(n)$  is nil or  $n'$  lies betw  $\text{prev}(n)$  &  $n$ :  
 its predecessor  $\text{prev}(n) = n'$   
 Successor when notified

## c) Stabilize()

→ Sets up the successor

→ It is run by every node periodically

$$x = \text{prev}(\text{next}(n))$$

If  $n$  in  $(n, \text{next}(n))$  then  $\text{next}(n) = x$

notify  $\text{next}(n)$

→ also each node  $n$  periodically updates a random finger entry

Pick a random  $i$  in  $[0, m-1]$

Lookup successor  $(n + 2^i)$

## Lec - MapReduce

# How to run large scale distributed:

Computations over key-value stores

⇒ Map Reduce Programming Abstraction

Key value stores in Cloud Computing → is delivery of computing resources - like servers, databases,

→ models: SaaS, PaaS, IaaS networking, software &

deployment models: public cloud

private "

Community "

hybrid "

Analytics

- over the

internet on

demand

Usually on a

pay as you

go basis!

→ Characteristics: On demand

Scalable

- (dynamic acquiring - elastic computing (and storage of compute resources) service)

→ What is a cloud?

Cloud = Lots of Storage +  
↓  
Services include,  
Compute cycles nearby

(i) managed clusters for distributed computing

(ii) managed distributed datastores

3 Tier architecture

(i) frontend

(ii) Compute nodes

(iii) Storage nodes

→ How to do processing on large datasets?  
↓

• distribute over multiple machines

• parallel computation  
↓

Map Reduce = it is an abstraction (a model) for doing large-scale computation easily in the cloud

Input : a set of key/value pairs

User supplies two functions:

(i)  $\text{map}(K, V) \rightarrow \text{list}(K_1, V_1)$

(ii)  $\text{reduce}(K_1, \text{list}(V_1)) \rightarrow V_2$

Intermediate

Key/Value

pairs

Output is the set of  
( $K_1, V_2$ ) pairs

Example:

```
map(key, value):  
    // Key : metadata record #  
    // Value : (URL, size, ---)
```

for each (URL, size) in size:  
 emit(URL, size)

reduce (key, values):

```
// Key : target  
// Values: list of pages that link to it  
result = concatenate(values)  
emit(key, result)
```

### Resource

- NOTE: Reducer Manager (assigns map and reduce tasks to servers)

Example: Yarn, Kubernetes

→ Treats each server as a collection of  
Containers

fixed CPU and memory (ex. Docker)

→ has 3 main components

(i) Global Resource Manager (RM)

→ cluster scheduling



(ii) Per server Node manager (NM)

→ daemon and server specific  
functions

- (iii) Per application (job) Application Master (AM)  
→ container negotiation with RM and NMs  
→ handling task failures of that job

- Fault Tolerance

(i) If NM fails ← (NM heartbeats RM)  
→ it stops sending heartbeats to RM  
→ RM informs AM

(ii) If AM fails ← (AM heartbeats to RM)  
→ RM restarts AM

(iii) If RM fails → use old checkpoints and bring up secondary RM

- Slow Servers

→ aka stragglers  
→ delay the entire job as all maps must finish before Reduce tasks can start

→ Reason: (i) bad disk

(ii) low memory

(iii) low CPU

(iv) poor network

(v) temporary Overload

→ How to handle?

- Uses speculative execution (MapReduce)

If one task is too slow:

(i) run another copy of the same task on a different machine

ii) Whichever copy finishes first is accepted

iii) The other copy is killed

- Task Scheduling

- try to schedule in the server where it is available

- If not possible, in same rack  
If not " , anywhere

- ⇒ otherwise, just delay !!

## L - Peer to Peer

- P2P Systems

- refer to applications that take advantage of resources (Storage, cycles, content, --)  
available at the end systems of the Internet

- Overlay networks

- refer to networks that are constructed on top of another network (eg. IP)

- P2P overlay network

- any overlay network that is constructed by the Internet peers in the application layer on top of the IP network

- Properties : (i) self organizing

- (ii) load balancing

- (iii) dealing with free riders

## # Lookup Issue

## # Types of overlay networks

(i) Structured overlay networks — DHT

(ii) Unstructured overlay networks

(iii) Overlay multicast networks

Join network with  
loose rules

Structured V/s Unstructured

(i) Specific structure

(i) no specific structure

— for data/placement

(ii) insertion/deletion have some overhead

(ii) no overhead in terms of both

(iii) fast lookup

(iii) lookup has high overhead

(iv) Complex queries

difficult to support

(iv) Complex queries

Supported

like keyword,

Range &amp; attribute

° done through queries

flooding or random walks

⇒ example: centralized

directory based P2P

Systems, Pure P2P Systems,

hybrid P2P Systems

## # Centralized Directory based P2P Systems

a) all peers connected to central entity



→ Central entity Server as:

(i) to provide the service

(ii) Some kind of Index/group database

(iii) lookup/routing table

ex. Napster, BitTorent



→ Single point of failure

limited scalability (server farms with

load balancing)

→ Query is fast  $\Rightarrow$  upper bound for duration

can be given

b) Gnutella

→ pure P2P

→ no routing intelligence

→ Constrained broadcast

(i) TTL typically set to 7

ii) packets have unique ids to detect loops

NOTE: fully distributed

(bandwidth  $\rightarrow$  edge: not a physical link  
↑ wasted)

Scalability:  
limited

Scope flooding

$\rightarrow$  Query flooding - fwd Query msg  
Query bit sent

Over reverse path

Flooding  $\xrightarrow{(i)}$  Random Walk

$\rightarrow$  scale free graph

- a network where

(i) a few nodes have very high degree } follows  
(many connections)

(ii) most nodes have very low degree } power  
law distribution

$\rightarrow$  Random walk

(i) Choose a node/neighbor at random

(ii) Move to it

(iii) Repeat

(iv) avoid immediately going back to the  
last node (to prevent loops)

Spreads the search through the  
network without flooding

$\rightarrow$  Degree biased Random walk

(i) pick the neighbour with the highest  
degree that hasn't been visited yet

### (ii) Replication

$\rightarrow$  Spread copies of objects to peers : more  
popular objects can be found easier

$\rightarrow$  replication strategies

a) Owner replication

b) Path II

c) Random II

## (iii) Random Walkers

- don't send to all neighbours
- only to a subset of its neighbours

## (IV) Unstructured Informed Searches

## a) informed networks

- nodes do smarter routing by using info they already learned

## • 2 types

→ Local indices

a) each node keeps info abt files within r hops

b) When querying, it asks neighbours just outside r

c) This reduces flooding because we search only a local zone

→ Routing Indices

a) Each node keeps an index of

◦ what kinds of files each neighbor likely has

◦ how many matching files each neighbor has for a given theme

b) each neighbor gets a goodness score

c) The node forwards the query to the best neighbour

## (V) Free Riding — Problem in P2P networks

→ Free riders are nodes that

- download but do not share
- or share useless/empty data

## # P2P Systems

Structured P2P Systems

even rare obj are found

efficiently!

## a) DHT

→ Hash-table functionality in a P2P network: lookup of data indexed by keys

ex. rooted queries done by Chord, Pastry etc

## b) Routing Challenges

→ define a useful key nearness metric

→ keep the hop count small

→ 11, 11, 11 tables

→ stay robust despite rapid change

## Freenet — Emphasizes anonymity

## Chord — " efficiency and simplicity

→ accounts hop count + network locality

## c) Pastry

◦ for routing efficiency

→ self organizing overlay network of nodes

→ Pastry takes into account network locality

→ eg IP routing hops

→ no ring like structure ◦ routing based on numeric closeness of identifiers

→ Identifier space.

- $(0 \rightarrow 2^m - 1)$  : m bit address

- but each digit in ID is  $2^b$  base

- keys are stored on the node whose ID is numerically closest to the key

$$\text{distance} = |\text{nodeID} - \text{keyID}|$$

→ Routing Graph

- takes  $\lceil \log_2(N) \rceil$  steps!! (for finding the node with the key)

- each node has 3 main elements

- (i) routing table (R)

- $M \left( \frac{\log_2(N)}{b} \right)$

- similar to Chord's finger table

- stores links to id-space

rows

- $2^b$  entries (ii) leafset

~~(per row!!)~~

- stores nodes close in id-space

- i.e.  $2^b + 1$  columns (iii) neighbourhood set

~~(to represent the digits)~~

- nodes that are closer together in terms of network locality

- each row corresponds to a prefix length

- each col stores nodes with the same prefix but differing next digit → i.e. prefix length == row no then column should differ in row no index

~~(Only 1 node exists !! per rows)~~

→ Leaf Set L

- It keeps  $|L|$  nodes numerically closest to  $n(|L|/2)$   
Smaller and  $|L|/2$  larger)

- imp info for routing : acts as chord's successor list in case of node failure

→ Neighbourhood Set (M)

- concerned with nodes that are close to the current node with regard to the network proximity metric

i.e how close they are in real network

## Routing Rules → When node A receive msg for key D

- a) If D is in A's leaf set range → Use the leaf set
  - b) else → use routing table
  - Compute  $L = \text{shared prefix length } (A, D)$
  - Then find to :
    - (i) routing table entry in row L
    - (ii) column = digit of D at position L  
→ if it exists
  - Else
    - (i) Choose any node from  $(R \cup TM \cup L)$
    - (ii) with equal or longer prefix match
    - (iii) and numerically closer to D than A
- ↓  
to prevent routing loops
- Prefix!
- \* better  
 as query routed to  
 to a node with a longer prefix than our Node  
 or to  
 a numerically closer with some

## → Routing Performance

### Case 1: R-T

- search space reduced by  $2^b$  factor at least
- destination reached in  $\log_2 N$  steps

### Case 2: L.S

- no. of hop + 1

### Case 3: TULUR

- no. of hop + 1

⇒ Complexity of routing

$$= O(\log_2 N)$$

# Join and failure



failed leaf node

◦ find numerically

closet by routing

◦ build its 3 table based  
on the entries of nodes  
on the route !!

◦ contact a leaf

node on the side

◦ of the failed

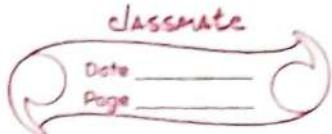
node and add

◦ appropriate new  
neighbour

failed table entry

◦ contact a live entry with same prefix  
as failed one

◦ & keep trying with longer prefix



## # Self Organisation

- o A new node  $i^2$  knows another node  $K$  as closer to it via proximity

$\Rightarrow$  Consider it's neighbouring set as their own!

- o thus via  $K$  it finds a numerically closer node using routing  $\xrightarrow{\text{with key} = n}$  takes it's leaf set as it's own
- o next, for R-T it takes routing info from the nodes visited in that routing process!

## # Lazy Detection of Node Failures



$\rightarrow$  as failure detected when a communication attempt with another node fails

↓  
and since routing involved  $\Rightarrow$  lazy detection

## # Node Departure

- asks when node departs:

(i) remaining nodes need to repair their neighbourhood set  $M$

(ii) they do this by asking other nodes for their  $M$

## # Locality

$\rightarrow$  It means that the routing algo tries to use nodes that are physically close in the real network, not just close in the id space!

# Locality in Routing Table

## New nodes copy routing tables from nearby nodes →  
## good entries propagate

# Route Locality

→ Each hop gets closer to target & is physically  
near → no loops, efficient path

# Nearest among k nodes

→ Switch to numeric closeness when near target  
to pick physically closest replica

### Extra Notes

#### a) GFS (Google File System)

- DFS designed by Google — for large scale data processing
- optimized for huge files, high throughput and fault tolerance
- Uses a single master service that stores metadata
- files — Split 64 MB chunks stored on chunkservers
- each chunk replicated (usually 3 times) for fault tolerance

- designed for write once, read many workloads
- ensures reliability through heartbeat messages, chunk versioning and automatic re-replication

## b) Andrew File System

- DFS developed at CMU widely used in universities
- focuses on scalability and efficient file sharing across many clients
- uses whole file caching on clients → entire file copied to client machine
- reduces server load because most operations happen on cached copy
- uses callbacks: server tells clients when cached file is invalid → strong consistency
- file space organized into volumes, which can be moved easily for load balancing
- Authentication via kerberos, supports secure access

- GHS

- distributed version of Prim's algo
- nodes form ~~frags~~ fragments & repeatedly find their least-weight-outgoing-edge (LWOE)
- fragments absorb or merge based on their levels, eventually forming the MST

- each fragment has a level

- 2 fragments of same level merge

$$\rightarrow \text{new level} = L+1$$

- lower<sup>level</sup> fragment absorbed into higher level fragment

- level tracks fragment size & controls merging

- Purpose of test/accept/reject messages

To decide if an edge is an outgoing edge

a) test - Is this edge connecting to a different fragment?

b) accept - yes, outgoing

c) reject - not, internal edge

- Uses broadcast & convergence

To ensure allow all nodes in a fragment to

report their ~~their~~ candidate edges to the root

so the fragment can find its LWOE

- branch and rejected edges



part of edges

↳ Internal edge that cannot  
be used for MST (once  
rejected, permanent)

o Replication

→ Why? availability

reliability

performance

when functioning even nodes fail

→ diff bet linearizability and sequential consistency

realtime order

is preserved

Strongest practical

model

only program

order preserved

real time may

differ

Linearizability  $\Rightarrow$  Sequential



→ eventual consistency

- a weak model where replicas may be temporarily inconsistent, but eventually converge if no new updates occur

- used in large scale systems

- example: gossip, Dynamo

→ Passive (primary backup)

Primary executes operations, sends updates to backups;  
implements linearizability

→ Active

All replicas execute operations in same order via  
total order multiset; gives sequential consistency !!

→ Client centric consistency model

- consistency guarantees from the perspective of a single client

- read after write

- monotonic read

- monotonic write

- write follow read

- Self Stabilization

- Self stabilizing system

- Start from arbitrary state

- Converges to a legal state in finite time

- ε, stays there

- 2 properties

- (i) Closure: from a legal state, system stays legal

- (ii) convergence: from any state, system eventually reaches a legal state

- Dijkstra's ~~token~~ ring <sup>token</sup> is a self stabilizing system

- Nodes arranged in a ring with status

- exactly one privileged node should exist

- (the token)

- who ensure that regardless of initial states,

- eventually exactly one token exists ε

- Circulates — achieving self stabilization

→ Why an exceptional node in Dijkstra's token ring?

If all nodes identical, Symmetry prevents system from ensuring a unique privileged node  
one node must behave differently to break Symmetry

→ What type of faults handled by self-stabilization

Transient faults (like memory corruption) not permanent failures

- Distributed File Systems

→ a remote file storage system that allows clients to access files transparently over a network as if they were local

→ transparency requirements satisfied

(i) access transparency

(ii) location "

(iii) migration "

(iv) replication "

(v) performance & scaling

→ NFS - Client-Server model

RPC based

stateless in V3

stateful in V4

GIPS - cluster

based DFS

→ Caching improves performance by reducing remote access

Introduces consistency problems since multi-cached copies may differ

→ file striping — (PVFS) (object based)

Striping a file into multiple blocks  
and storing them across multiple servers  
to establish parallel access