

Deadlock handling Strategies \Rightarrow

- * Prevention - inefficient
 - * Avoidance - impractical
 - * Detection & Resolution
 - ① Progress (no undetected deadlocks)
 - detected and resolved
 - ② Safety (No false deadlocks) - No phantom deadlocks

} essential conditions for a deadlock detection algorithm

Models of deadlocks \Rightarrow A process might require a single resource or a combination of resources (for its execution)

D) The Single resource model — A process can have one outstanding request for only unit of one resource.

g) AND Model - A process may request more than one resources simultaneously and the request is satisfied only after all the requested resources are granted to the process.

3) OR Model - A process may request more than one resources simultaneously and the request is satisfied if any one of the resources requested is granted.

- A cycle doesn't imply deadlock in OR Model unlike AND Model
- Presence of
→ A Knot indicates a deadlock
- In a WFG, a vertex v is a Knot, if for all u : u is
reachable from v : v is unreachable from u
- No path originating from a Knot have deadends.

1. Each of the processes in the set S is blocked
 2. The dependent set for each process in S is a subset of S
 3. No grant message is in transit between any two processes in the set S .
(a process may release a resource)

4) AND-OR Model - In this model, we need $x_1, y_1, z \dots$

5) $\binom{P}{q}$ Model \Rightarrow requires p resources from a pool of q resources.

Request to obtain any p available resources from a pool of q resources.

Special Case:- $\binom{P}{p}$ Model - corresponds to AND Model

$\binom{1}{P}$ Model - corresponds to OR Model.

6) Unrestricted Model - No assumptions are made regarding the underlying structure of resource requests.

detect a deadlock based on the global state.

Deadlock Detection Algorithm:-

- * Path pushing - arrive at a global WFG.
(all requirements for all the processes)
- * edge-chasing - probe message is sent in the opposite direction of request messages along the edges to detect cycle.
- * Diffusion computation - query messages to the adjacent processes; it'll reply if it is in deadlock
 \rightarrow deadlock detection is diffused.
- * global state detection - arrive at global WFG and analyse for cycles or knots.

Path pushing algorithms -

- maintain an explicit global WFG.
- each process will send WFG through one process.

Edge-chasing algorithms -

- presence of a cycle in a distributed graph structure is verified by propagating special messages called probes along the edges of the graph.
- if the probe message comes back, then there's a cycle.

Diffusion-Computation-based algorithm -

- Deadlock detection computation is diffused through the WFG of the system, makes use of echo algorithm to detect deadlocks.
(Reply is an echo)

Global state detection based algorithms -

- Taking a snapshot of the system and examining it for the condition of a deadlock (cycles and knots)

* Mitchell and Merritt's algorithm for a single resource model \Rightarrow Edge chasing algorithm.

- Probes are sent in the opposite direction to the edges of WFG.

1) Only one process in a cycle detects the deadlock - This can be improved by including priorities and the lowest priority process can be aborted.

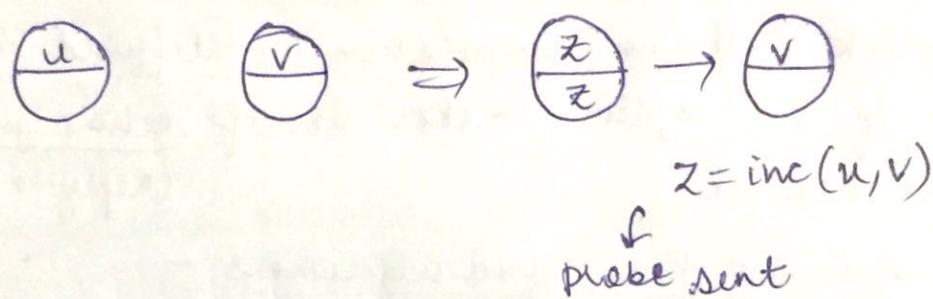
2) In this algorithm, a process that is detected in a deadlock is aborted spontaneously, even though under this assumption, phantom deadlocks cannot be excluded.

In the absence of spontaneous aborts, genuine deadlocks will be detected.

Each node of WFG, two local variables = labels
↓
a public label ↓
a private label.

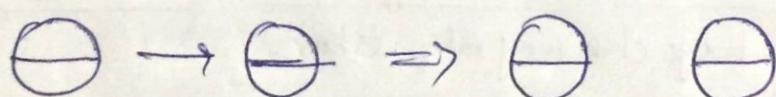
- Block creates an edge in WFG - Two messages are needed, one resource request and one message back to the blocked process to inform of the public label of the process it is waiting for

Block :-



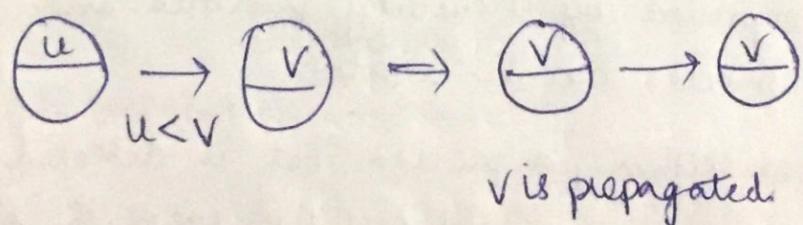
- Activate denotes that a process has acquired the resource from the process it was waiting for.

Activate :-



- Transmit propagates larger labels in the opposite direction to the edges by sending a probe message.

Transmit :-



Detact means that the probe with the private label of some process has returned to it, indicating a deadlock

Detact : $\frac{u}{u} \xrightarrow{p} \frac{u}{u} \Rightarrow \frac{z}{z} \xrightarrow{p'} \frac{v}{v}$

Initially we have $u=v$ for all processes

The only request that change $u & v$, Block and activate.

$\rightarrow u$ has been propagated from p to p' .

\rightarrow probe transmitted only by a deadlocked process

Chandy - Misra - Haas algorithm for the AND Model \Rightarrow

- * Edge-chasing algorithm
- * Probe is sent out in the form of triplet (i, j, k)
[it is sent along the edges]
[unlike in opposite direction]
- \downarrow
 $\begin{array}{l} \text{process } i \\ \text{initiating} \end{array} \quad \begin{array}{l} \text{probe is sent from} \\ \text{process } j \text{ to } k \text{ across} \\ \text{sites.} \end{array}$
- \rightarrow A deadlock detection initiated for p_i and is being sent by the home site of process p_j to the home site of process p_k .
- \rightarrow A probe message travels along the edges of the global WFG and a deadlock is detected when a probe message returns to the process that initiated it.
- \rightarrow A process p_j is said to be dependent on another process p_k if there exists a sequence of processes.
$$P_j, P_{j1}, P_{j2}, \dots, P_{jm}, P_k$$

could be null also
(there might be direct dependency also)

Each process p_i maintains a boolean array dependent_i , where $\text{dependent}_i(j)$ is true only if p_i knows that p_j is dependent on it.

Initially, $\text{dependent}_i(j)$ is false for $\forall i, j$,

If p_i is locally dependent on itself, then declare a deadlock
else for all p_j and p_k such that

- (a) p_i is locally dependent on p_j and
- (b) p_j is waiting on p_k
- (c) p_j and p_k are on different sites.

Send a probe (i, j, k) to the home site of p_k .

On the receipt of a probe (i, j, k) , the site takes the following actions -

if

- (d) p_i is blocked and
- (e) $\text{dependent}_k(i)$ is false.
- (f) p_k has not replaced p_j updating p_j

then

begin

$\text{dependent}_k(i)$ is true

if $k = i$ then declare that p_i is deadlocked
else for all p_m and p_n such that

- (a) p_k is locally dependent on p_m and
- (b) p_m is waiting on p_n
- (c) p_m and p_n are on different sites

Send a probe (m, n) to the home site of p_n

end

- should not require extra additional communication channel.

Chandy-Misra-Haas algorithm for the OR Model \Rightarrow

\rightarrow Diffusion computation

\rightarrow 2 types of messages \Rightarrow query (ijj, k) and

$p_j \downarrow$
 p_j sends a query to p_k

reply (ijj, k)

$p_j \downarrow$
 p_j sends a reply to p_k

\rightarrow If an active process receives a query, it ignores the message.

Process p_k maintains a boolean variable $\text{wait}_k(i)$ that denotes the fact that it has been continuously blocked since it received the last engaging query from process p_i .

Each process has a local variable $\text{num}_k(i)$

\rightarrow Initiate a diffusion computation for a blocked process p_i

\rightarrow Send query (i, i, j) to all processes p_j in the dependent set DS_i of p_i

$$\text{num}_j(i) = |DS_i|, \quad \text{wait}_j(i) = \text{true}.$$

\Rightarrow When a blocked process p_k receives a query (ijj, k) , if this is the engaging query for process p_i

then send query (i, k, m) to all p_m in its dependent set $|DS_k|$

$$\text{num}_k(i) = |DS_k|; \quad \text{wait}_k(i) = \text{true}$$

else if $\text{wait}_k(i)$, then send a reply (i, k, j) to p_j

When a process P_k receives a reply ($i j k$)

→ if $\text{wait}_k(i)$ then $\text{num}_k(i) = \text{num}_k(i) - 1$

if $\text{num}_k(i) = 0$, then if $i = K$, then

"Declare a deadlock"

else

send reply ($i k m$) to the process P_m (which sent the engaging query)

* Messages are different from the underlying computation messages and the message communication is also transparent.

→ If every member of the dependent set sends a reply, then "a deadlock is detected".

Termination Detection → not easy to determine.

Halting problem in Turing Machines.

Two distributed computations

① the underlying computation — basic messages

② termination detection algorithm — control messages

* A termination detection (TD) algorithm must ensure the following —

① Execution of TD algorithm cannot independently indefinitely delay the underlying computation, must not

freeze the underlying computation.

② The TD algorithm must not require addition of new communication channels between processes.

⇒ A distributed computation has the following characteristics-

① Processes can be only one of two states - active and idle
or (busy and passive)

② An active process can become idle at any time

③ An idle process can become active only on the receipt of a message

④ Only active processes can send messages

⑤ A message can be received by a process when the process is either of the two states : active or idle

⑥ The sending of a message and receipt of message occur as atomic access. (not affected by any other external means)

$$* (\forall i :: P_i(t_0) = \text{idle}) \wedge (\forall i j :: C_{ij}(t_0) = 0)$$

Termination Detection using Distributed Snapshots -

* Global Snapshot [collection of individual snapshots]
↓
Distributed.

The algorithm is defined by the following four rules-
Each process i applies one of the rules when it is applicable.

R1: When process i is active, it may send a basic message to process j at any time by doing -

Send a $B(x)$ to j

↓
basic message sent by the process.

R2: upon receiving a $B(x')$, process i does -

$$\text{let } x = x' + 1$$

if ($i \Rightarrow \text{idle}$) \rightarrow go active [// become active on receiving a message]

R3: when a process i goes idle, it does \Rightarrow

$$\text{Let } x = x + 1$$

$$\text{Let } k = i$$

Send message $\underbrace{R(x, k)}_{\substack{\text{control} \\ \text{message}}}$ to all other processes.

[with timestamp
process no.]

Take a local snapshot for the request of $R(x, k)$.

R4: Upon receiving message $R(x', k')$, process i does \rightarrow

$[(x', k') > (x, k) \wedge (i \text{ is idle}) \rightarrow$

$$\text{let } (x, k) = (x', k')$$

Take a local snapshot for the request by $R(x', k')$

\square requesting process is earlier to the current process.

$((x', k')) \leq (x, k) \wedge (i \text{ is idle}) \rightarrow \text{do nothing}$

\square

$(i \text{ is active}) \rightarrow \text{let } x = \max(x, x')$

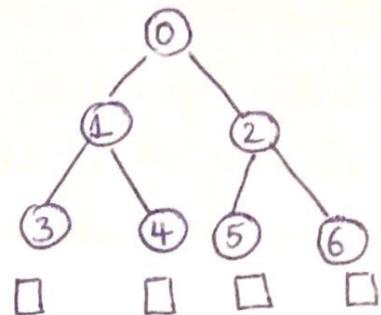
Just a tree datastructure
+
colouring scheme
(black + white)
tokens

- associated with leaf nodes
- sends a message to its parent which means one of the processes hasn't terminated.

- before termination - sends a black token to its parent

- An intermediate node root might receive one/more black tokens → infer termination not completed.

tokens
are passed
on once
all the
children nodes
are terminated.



- if the root node receives a black token, it repeats the whole process again
- This process terminates iff the root node receives a white token/s.

Spanning-tree-based termination-detection algorithm :-

- * Tokens - a contracting weave of signals that move inward from the leaves to the root.
- * Repeat signal - if a token weave fails to detect the termination, node P_0 initiates another round of termination detection by sending a signal called "Repeat" signal to the leaves.
- * The nodes which have one or more tokens at any instant form a set S
- * A node j is said to be outside of set S if j does not belong to S and the path (in the trees) from the root to j contains an element of S
- Note that all nodes outside of S are idle

Algorithm -

- ① Initially each leaf process is provided with a token. The set S is used for book keeping tokens which processes have the token.
- ② Initially, all processes and tokens are white.
- ③ When a leaf node terminates, it sends the token it holds to its parent process.
- ④ A parent process will collect the token sent by each of its children. After it has received a token from all of its children and after ~~which~~ it has terminated the process sends a token to its parent.
- ⑤ A process turns black, when it sends a message to some other process. When a process terminates, if it is black, it sends a black token to its parent.
- ⑥ A black process turns back to white after it has sent a black token to its parent.
- ⑦ A parent process holding a black token (from one of its children) sends ~~only~~ a black token to the parent node to indicate that a message passing was involved in the subtree.
- ⑧ Tokens are propagated to the root in this fashion. The root after receiving a black token, will know that a process in the tree had sent a message to some other process. Hence, it restarts the algorithm by sending a repeat signal to all of its children.
- ⑨ Each child of the root propagates the repeat signal to

each of its children and so on until the signal reaches the leaves

- (i) The leaf nodes restart the algorithm on using the "Repeat" signal.
- (ii) The root concludes the termination has occurred - if
 - a) If it is white. (Should satisfy all the conditions)
 - b) If it is idle.
 - c) If it has received a white token from all of its children

[Sukumar Ghosh]

Termination detection for diffusion computation \Rightarrow

Environment node, internal nodes
 \rightarrow signals, acks.

\Rightarrow For an edge (i, j) , the difference between the number of signals sent by i and the number of acks received from j will be called a deficit.

\Rightarrow A process keeps track of two different types of deficit

$C =$ total deficit along its incoming edges

$D =$ total deficit along its outgoing edges

By definition, these deficits are nonnegative integers.

INV1: $(C \geq 0) \wedge (D \geq 0)$

Environment Node $\Rightarrow C = 0, D = k$

For every other node in the system, the proposed signalling structure scheme processes the following invariant.

$$\underline{\text{Inv 2}}: (c > 0) \vee (d = 0)$$

The sending of an ack, reduces the senders deficit c by 1

$$(c - 1 \geq 0) \wedge (c - 1 > 0 \vee d = 0)$$

$$= (c \geq 1) \wedge (c > 1 \vee d = 0)$$

$$= (c > 1) \vee (c \geq 1 \wedge d = 0)$$

$$= (c > 1) \vee (c = 1 \wedge d = 0)$$

Algorithm :-

- Sukumar Ghosh

α (some 1980 paper)

Program detect {for an internal node i }

define c, d : integer

m : (signal, ack) {represents the type of message received}
state: (active, passive)

initially $c = 0, d = 0, \text{parent}(i) = i^*$

do ($m = \text{signal}$) $\wedge (c = 0) \rightarrow c = 1, \text{state} = \text{active},$
parent = sender
[] $m = \text{ack} \rightarrow D := D - 1$

[] $(c = 1 \wedge d = 0) \wedge (\text{state} = \text{passive}) \rightarrow \text{send ack to sender}$.

[] $(m = \text{signal}) \wedge (c = 1) \rightarrow \text{send ack to the sender}$

The environment node initiates the termination detection algorithm by making $D = k$ ($k > 0$) and detects termination by the condition $D = 0$ after it receives k acks.

Leader Election →

Leader → focus of control = system wide management

It is tempting to compare leader election with the problem of mutual exclusion

* When a leader fails, you have to elect a new leader.

Comparison :- (ME vs Leader election)

1. Failure is not an inherent part of mutual exclusion algorithm.
2. Starvation is an irrelevant issue in leader election
[A leader, which ceases to fail, remains a leader...]
but that leads to starvation in the case of ME
3. If leader election is viewed from the perspective of mutual exclusion, the exit from the critical section is unnecessary.

[there's no point of leaving the CS in Leader election]

Let $G_1 = (V, E)$, each process $i \in V$ has a unique identifier
→ $L(i)$ identifies the leader.

→ $ok(i)$ indicates whether process is active/not faulty.

$$1. \forall i, j \in V : ok(i) \wedge ok(j) \Rightarrow L(i) = L(j)$$

$$2. L(i) \in V \quad (\text{Leader also belongs to the } \cancel{\text{one}} \text{ set}) \\ \text{of vertices.}$$

$$3. ok(L(i)) = \text{true}$$

Bully Algorithm →

Assumptions

- ① Communication links are fault-free
- ② Process can only fail by stopping

③ Failures can be correctly detected using some mechanism like timeout.

Messages: Election, reply and leader.

Algorithm:

Program $b_{i,V}$ { program for process $i \in V$ }

define failed: Boolean { set if the failure of the leader is detected }

L : process { identifies the leader }

m : message { election | leader | reply }

state: idle | wait for reply | wait for leader.

initially $\forall i \in V, state = idle, failed = false$

do (failure of $L(i)$ is detected) \wedge (failed = false) \rightarrow
failed = true,

[] failed = true $\rightarrow \forall j > i ;$ send election to $j,$
 $state = wait for reply,$
 $failed = false;$

[] state = idle $\wedge (m = election) \rightarrow$ send reply to sender,
 $failed = true;$

[] ($state = wait for reply$) $\wedge (m = reply) \rightarrow state = wait for leader$

[] ($state = wait for reply$) \wedge timeout $\rightarrow L(i) = i, state = idle,$
send leader(message) to all.

[] ($m = leader$) $\rightarrow L(i) = sender, state = idle$

$[] (\text{state} = \text{wait for leader}) \wedge \text{timeout} \rightarrow \text{failed} = \text{true}$,
 $\text{state} = \text{idle}$.

or

Chang-Roberts Algorithm \Rightarrow Leader election algorithm for a unidirectional ring.

define token : process id, color $\in \{\text{red, black}\}$

Initially all processes are red and i sends a token $\langle i \rangle$ to its neighbor

(for a red process i)

do token $\langle j \rangle$ received $\wedge (j < i) \rightarrow$ skip $\{j\}'s$ token removed,
so. j quits }

$[] \text{token } \langle j \rangle \text{ received } \wedge (j > i) \rightarrow \text{send } \langle j \rangle,$
color = black $\{i$ quits $\}$

$[] \text{token } \langle j \rangle \text{ received } \wedge (j = i) \rightarrow L(i) = i \{i$ becomes
the leader $\}$

or

(for a black process)

do token $\langle j \rangle$ received \rightarrow send $\langle j \rangle$

or

Consensus and agreement algorithm \Rightarrow

* commit decision in a database system, commit/reboot

\Rightarrow Computers controlling a spaceship's engine - proceed/abort.

Feasibility of designing algorithm to reach agreement under various system models and failure models. →

Assumptions underlying the study of agreement algorithms

Failure models → among the n processes in the system at most of processes can be faulty

The various failure models — fail-stop send omission and receive omission and Byzantine failure.

* Consensus and Agreement algorithms →

Assumptions underlying the study of agreement algorithms :-

* Failure Models — n processes f -processes

↓
Failstop, send omission, Receive omission and Byzantine failure.

* Synchronous / Asynchronous communication —

→ Impossibility of reaching agreement in asynchronous system in any failure model

→ In a synchronous system, the scenario in which a message has not been sent can be recognised by the intended recipient.
it'll know.

* Network Connectivity —

The system has full logical connectivity

* Sender identification — A process that receives a message always knows the identity of a sender message.

* channel reliability — The channels are reliable and only the processes may fail

(under one of the various failure models).

→ Two different kinds of messages — authenticated & non-authenticated messages.

Authenticated Vs Non-authenticated messages —

→ [we'll be dealing with only non-authenticated messages]

Messages with unauthenticated messages, when a healthy process relays messages to the processes

- (i) it can forge the message and claim that it was received from another process and
- (ii) it can tamper with the contents of a received message and before relaying it, ~~whereas~~ when a process receives a message, it has no way to verify its authenticity.

Unauthenticated message is referred to as — Oral / unsigned message.

Agreement Variable — may be boolean or multivalued and need not be an integer.

→ Byzantine agreement and other problems —

* Byzantine agreement problem → Requires a designated process called the source process with an initial value to reach an agreement with the other processes about its initial value, subject to the following conditions.

Agreement - All non faulty processes must agree on the same value.

Validity - If the same process is nonfaulty, then the agreed upon value by all the nonfaulty processes must be the same as the initial value of the source.

Termination - Each nonfaulty process must eventually decide on a value.

* Concensus problem \rightarrow Each process should have an initial value and all the correct processes must agree^{to} a single value

Agreement \rightarrow All nonfaulty processes must agree on the same (single) value.

Validity \rightarrow If all the nonfaulty processes have the same initial value, then the agreed upon value by all the nonfaulty processes must be that same value.

Termination \rightarrow Each non faulty process must eventually decide on a value.

* The interactive consenting problem - Each process has an initial value and they must agree upon a set of values with one value for each process.

Agreement \Rightarrow All the non faulty processes must agree on the same array of values $A[v_1, v_2 \dots v_n]$

Validity \Rightarrow If process i is non faulty and its initial value is v_i , then all the non faulty processes

agree on v_i as the i th element of the array A .
If process j is faulty, then the non-faulty processes
agree on any value v_i for $A[j]$.

Termination — Each non-faulty process must eventually
decide on the array, A .

Equivalence of the problems and notation :-

→ The three problems defined above are equivalent
in the sense that a solution to any one of them,
can be used as a solution to the other two problems,
the equivalence can be shown using a reduction
of each problem to the other two problems.

→ If problem A is reduced to problem B, then a
solution to problem B can be used as a solution
to problem A in conjunction with the reduction.

If P is reducible to P' and P' is solvable so is P .

If P is reducible to P' and P is unsolvable, so is P' .

— Halting Problem (HP) of TM is unsolvable

If follows that every class of P' of yes/no problems,
if HP is reducible to P' , then P' is unsolvable.

= It is sometimes possible to device a solution to one
problem using a solution to another.

= Suppose that there exist solutions to consensus (C),
byzantine agreement (BG) and interactive consistency (IC)
as follows -

$C_i(v_1, v_2, \dots, v_n)$ returns the decision value of P_i in a run of the solution to consensus problem, where v_1, v_2, \dots, v_n are the values that the processes proposed.

$BG_j(j, v)$ returns the decision value of P_i in a run of the solution to the byzantine problem where j , the commander proposes the value v .
 $j \rightarrow$ commander.

$IC_i(v_1, v_2, \dots, v_n)[j]$ returns the j^{th} value of in the decision vector of P_i in a run of the solution to the interactive consistency problem where v_1, v_2, \dots, v_n are the values that the processes proposed.

IC from BG \rightarrow we construct a solution to IC from BG by running BG N times, once with each $P_i(i, j)$ acting as commander.

$$IC_i(v_1, v_2, \dots, v_N)[j] = BG_i(j, v_j) \quad (i, j = 1, 2, \dots, N)$$

C from IC \rightarrow For the case where a majority of processes are correct, we construct a solution to C from IC by running IC to produce a vector of values at each process, then applying an appropriate function on the vector's values to derive a single value.

$$C_i(v_1, \dots, v_N) = \text{majority} \left(IC_i(v_1, \dots, v_N)[1], \dots, IC_i(v_1, \dots, v_N)[N] \right)$$

$\forall i=1, 2, \dots, N$, where majority is defined as above.
 (majority over all values processes)

- BG from C - The commander j sends its proposed value v to itself and each of the remaining processes
- All processes run C with the values v_1, v_2, \dots, v_N that they receive (p_j may be faulty).

They derive,

$$BG_j(j, v) = C_i(v_1, \dots, v_N) \quad \forall i=1, 2, \dots, N$$

for each process.

A solution to one problem \Rightarrow Solution to other problems.

Agreement in a failure-free system (synchronous / asynchronous system)

In a failure-free system, consensus can be reached by collecting information from the different processes, agreeing at a decision and distributing the decision in the system.

In a synchronous system, in a constant number of rounds.

In an asynchronous system, in a constant number of steps

Agreement in (message passing) synchronous system with
failures \rightarrow

Concensus algorithm for crash failure

— kshemakalyani

↓
other processes need not be
knowing when a process crashes
or detected unlike fail-stop

\rightarrow n processes in the system, out of which f processes may crash ($f < N$)

Concensus variable, x is integer-valued

\Rightarrow (global constants)

integer $f \leftarrow //$ maximum no. of crash failures tolerated.

(local variables)

integer $x \leftarrow$ local variable

(1) Process P_i ($1 \leq i \leq N$) executes the concensus algorithm for upto f crash failures.

(a) for round from 1 to $f+1$ do

(b) if the constant value of x has not been broadcasted
then

(c) broadcast(x);

(d) $y_j \leftarrow$ value (if any) received from process j in
the round.

(e) $x \leftarrow \min_{y_j} (x_i y_j)_i$

I(f) output x as the consensus value.

The agreement condition is satisfied because in $f+1$ rounds, there must be at least one, in which no process failed.

~~consensus~~ round.

The validity condition is satisfied because processes don't send fictitious values in their failure model.

The termination condition is seen to be satisfied.

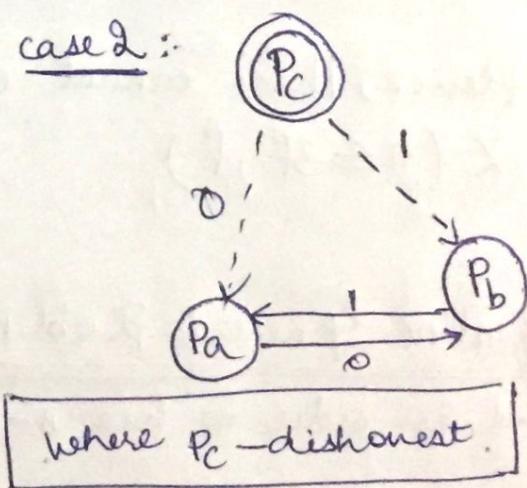
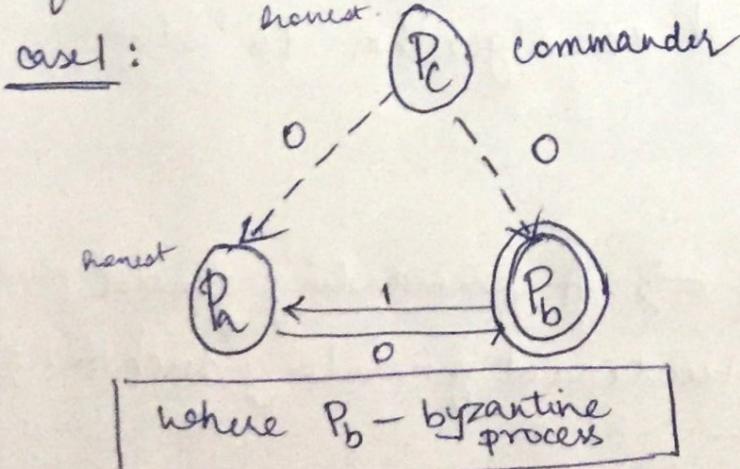
Consensus algorithm for Byzantine failures [Synchronous systems]

upperbound in Byzantine processes =

In a system of n processes, the byzantine agreement problem can be solved in a synchronous system only if the number of byzantine processes f is such that $f \leq \left\lfloor \frac{n-1}{3} \right\rfloor$

* with $n=3$ processes, the byzantine agreement problem cannot be solved if the number of byzantine processes,

$$f=1.$$



* P_a receives 0 from P_c & 1 from P_b P_a receives 0 from P_c & 1 from P_b
can't decide who is disloyal, because P_a gets the same numbers in both the cases.

$BG(3,1) \rightarrow$ not solvable, considering the two situations.



can be extended to n processes.



With n processes and $f > \frac{n}{3}$ processes, the Byzantine agreement cannot be solved, the correctness of this result can be shown using reduction.

Let $\mathcal{Z}(3,1)$ denote the Byzantine agreement problem for parameters $n=3$ & $f=1$, let $\mathcal{Z}(n \leq 3f, f)$ denote the Byzantine agreement problem for parameters $n(\leq 3f)$ and f .

A reduction from $\mathcal{Z}(3,1)$ to $\mathcal{Z}(n \leq 3f, f)$ need to be shown.

→ n processes divided into 3 groups, one group with all the faulty processes.

→ With this reduction in place, if there exists an algorithm to solve $\mathcal{Z}(n \leq 3f, f)$ i.e to satisfy validity, agreement and termination conditions, then there also exists an algorithm to solve $\mathcal{Z}(3,1)$ which has been seen to be unsolvable.

Hence, there cannot exist an algorithm to solve $\mathcal{Z}(n \leq 3f, f)$

Byzantine Generals Problem \Rightarrow A commanding general must send an order to his $n-1$ lieutenant generals, such that \Rightarrow

IC1: All loyal lieutenants obey the same order

IC2: If the commanding general is loyal, then

every loyal lieutenant obey the order he sends.

IC1, IC2 \rightarrow interactive consistency condition.

Concensus with oral messages \rightarrow The oral messages ~~wanted~~ model satisfies ~~for~~ the following conditions \Rightarrow

1. Messages are not corrupted in transit
2. Messages can be lost, but the absence of messages can be detected.
3. When a message is received (or its absence is detected) the receiver knows the identity of the sender (or the defaulter)

Let m be the number of traitors.

An algorithm that satisfies the IC1 and IC2 message model of communication in the presence of at most m traitors and be called an OM(m) algorithm.

The OM(0) algorithm is based on direct communication only.

\Rightarrow OM(m) Algorithm

In order to reach a concensus in the presence of m traitors, $n \geq 3m+1$ genuels are modelled. This algo is recursive \Rightarrow

$$OM(m) \rightarrow OM(m-1) \rightarrow OM(m-2)$$

Algorithm = OM(0) \rightarrow 0 traitors

1. Commander i sends out a value $v \in \{0,1\}$ to every lieutenant $j \neq i$.
2. Each lieutenant j accepts the value from i as the order from commander i .

Algorithm = OM(m)

1. Commander i sends out a value $v \in \{0,1\}$ to every lieutenant $j \neq i$.
2. If $m > 0$, then each lieutenant j after receiving a value from the commander starts a new phase broadcasting it to the remaining lieutenants using $OM(m-1)$. In this phase, j acts as a commander.

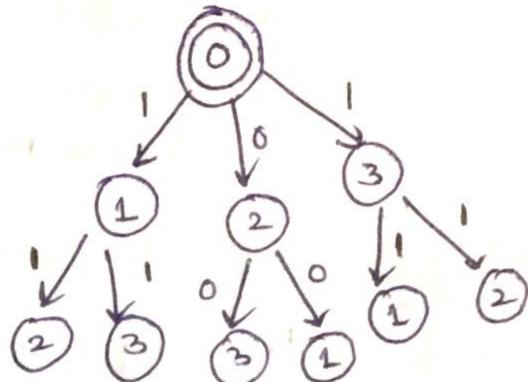
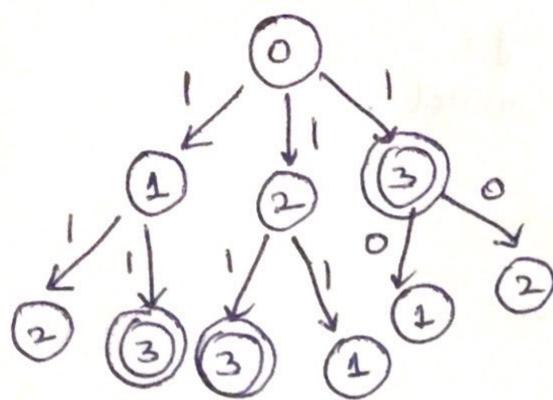
Each lieutenant thus receives $(n-1)$ values.

- Ⓐ a value directly received from the commander i of $OM(m)$
- Ⓑ $n-2$ values indirectly received from the $(n-2)$ lieutenants resulting from their broadcast of $OM(m-1)$

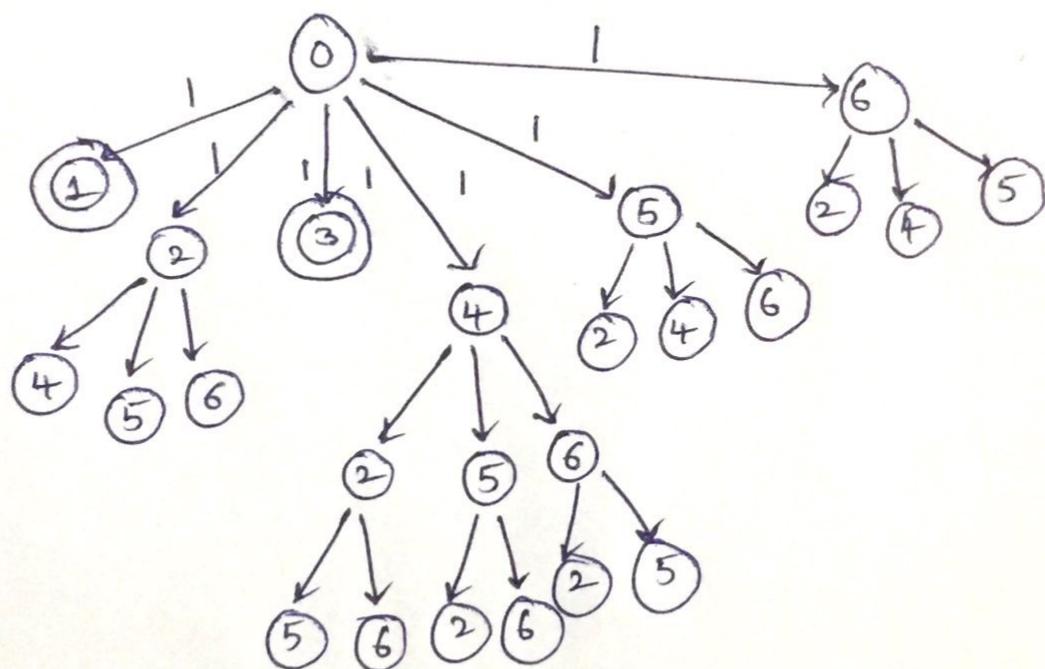
If a value is not received, then it is substituted by a default value.

3. Each lieutenant chooses the majority of the $(n-1)$ values received by it as the order from the commander.

OM(1)



OM(2)



Distributed Algorithms =

Algorithm → Input / No Input
→ Finiteness
→ Definiteness

Kepler graphic processors

- Global state
- Global time
- Non-determinism

Application — Telecommunication,
distributed information processing,
scientific computing,
real time process control.

Topology — undirected graph.

3 levels of topology abstraction.

— physical topology

— logical topology

— superimposed topology → for control algorithm execution
(ref to as a protocol)

- Distributed algorithm → all processes share computation and communication overhead.
- Centralised algorithm → one or few processes
For ex: client - server systems.

→ Symmetric and asymmetric algorithms



all the processes do similar kind of work.

— centralised algorithm

(few processes do majority of the work compared to other processes)

- Anonymous algorithm → An anonymous system in which process/professor code doesn't involve process numbers or processor nos. to make any execution decisions in the distributed algorithm.

→ structural elegance

→ hard to design

→ An anonymous algorithm is the algorithm which runs as an anonymous system and therefore doesn't use process identifiers or processor identifiers in the code.

Uniform Algorithm → A uniform algorithm does not use n , the number of processes in the system as parameters in its code.

↳ Scalability transparency

Adaptive Algorithm → Consider the context of problem x , in a system with n nodes, let k ($k \leq n$) be the number of nodes participating in the context of x . If the complexity of the algorithm can be expressed in terms of k rather than, in terms of n , the algorithm is adaptive.

Deterministic vs Non-Deterministic ↳ ^{Executions.}

Deterministic receive → process knows the sender process.

Non deterministic receive → process doesn't know the sender process.

Execution inhibition →

- Few protocols may affect the performance [even: event in sending delays or receive delays.]
- Blocking communicating primitives, freeze the local execution.
- Non blocking protocols, inhibitory or freezing protocols.
- Globally inhibitory
- Send/Receive inhibitory
- Internal event inhibitory

Synchronous and Asynchronous systems =

- Known upper bound on the message communication delay
- Known bounded drift rate for the clock
- Known upper bound to execute a logical step in the execution.
- ⇒ Asynchronous systems will not satisfy these three conditions.

* Online vs Offline algorithm — depends on data.

↓
unless entire data is available, execution doesn't happen
executed as data is generated.

* Failure models —

• t-fault tolerant — No. of components failed $\leq t$

↓
[processors or communication channels]

Process failure models — Fairstop — Processors may stop at a particular point, and this will be known to other processors also

- Crash failures.
- Receive omission
- Send omission
- General omission.

→ Byzantine failure or malicious failure with authentication

→ receive process knows the senders signature and can verify

Signature — cannot be manipulated

→ Byzantine or malicious failures.

↳ in the case of manipulated messages, without authentication
or cooked up - we cannot do anything

→ Communication failure models = crash failures
omission failures
byzantine failures.

* Wait

~~Wait-free algorithms~~ → each process has an upperbound

- resilient (continue to execute inspite of other processes)
- fault-tolerant systems.

* Communication channels — FIFO & NON FIFO .

* Complexity measures and metrics —

→ time and space complexity — lowerbounds (r, w)
upper bounds (O, o)
exact bound (δ)



→ message complexity

→ Space-complexity per node

→ system-wide space complexity.

→ Time complexity per node

→ system-wide time complexity.

→ Message complexity — number of messages,
size of messages,
message time complexity

Program structure - Hoare - CSP - concurrent sequential
proposed by problems process

$$* [G_1 \rightarrow CL_1 \parallel G_2 \rightarrow CL_2 \dots \parallel G_k \rightarrow CL_k]$$

guard = boolean value — if true, then executed
if all are false, then termination
takes place

Synchronous single-initiator spanning tree algorithm using flooding →

Synchronous BFS Spanning tree algorithm -

The code shown is for process P_i — $1 \leq i \leq n$
(local variables)

int visited, depth \leftarrow

int parent $\leftarrow \perp$

set of int neighbors \leftarrow set of neighbors.

(message type)

QUERY.

if $i = \text{root}$ then

visited $\leftarrow 1$

depth = 0

send query QUERY to Neighbors.

for round = 1 to diameter do

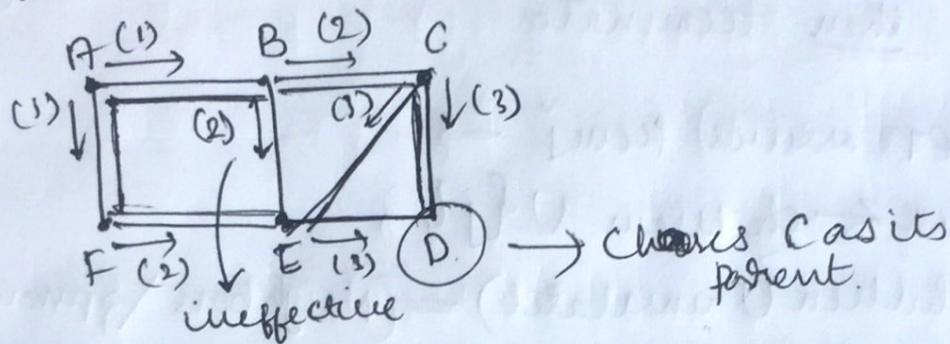
if visited = 0

if any QUERY message arrives then

parent \leftarrow randomly select a node from
which QUERY was received.

visited ← 1
 depth ← round
 send QUERY to neighbors \ {Senders of QUERY
 in the round}
 delete any QUERY message that is received in the
 round.

Example :-



→ Asynchronous single-initiate spanning tree algorithm →

Code for process P_i ; $1 \leq i \leq n$

(local variables)

int parent ← 1

Set of int children unrelated $\leftarrow \emptyset$

Set of int Neighbor \leftarrow set of neighbors.

(message type)

QUERY, ACCEPT, REJECT

→ When the predesignated root node wants to initiate the algorithm.

If ($i = \text{root}$ and $\text{parent} = 1$) then

Send QUERY to all neighbors.

parent = i

→ When QUERY arrives from j :

if parent = \perp then

parent = j

send ACCEPT to j

send QUERY to all neighbors except j

if (children \cup unrelated) = (neighbors \ {parent})

then terminate

↳ exclude parent

when ACCEPT arrives from j \rightarrow

children \leftarrow children $\cup \{j\}$

if (children \cup unrelated) = (Neighbors \ {parent})

then terminate.

when REJECT arrives from j \Rightarrow

unrelated = unrelated $\cup \{j\}$

if (children \cup unrelated) = (Neighbors \ {parent})

then terminate

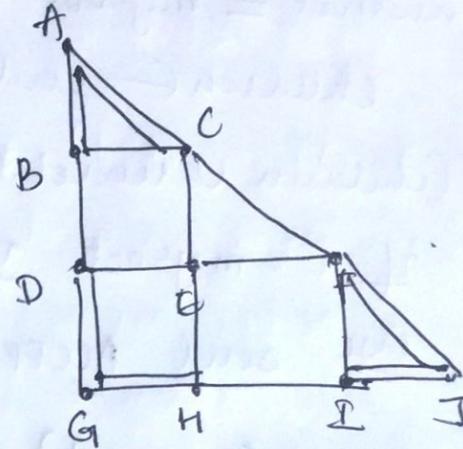
Asynchronous concurrent - initiator spanning tree algorithm using flooding

(local variables)

int parent, myroot $\leftarrow \perp$

set of int children_unrelated $\leftarrow \emptyset$

set of int Neighbors \leftarrow set of Neighbors



(message types)

QUERY, ACCEPT, REJECT

\Rightarrow when the node wants to initiate algorithm as a root
if (parent = \perp) then send QUERY(i) to all Neighbors.

parent, myroot $\leftarrow i$

\Rightarrow when QUERY(newroot) arrives from j

if myroot < newroot

then // discard . partial execution due to its low priority

parent $\leftarrow j$, myroot \leftarrow newroot

children_unrelated $\leftarrow \emptyset$

send QUERY(newroot) to all neighbors except j.

if neighbors = {j} then

send ACCEPT(newroot) to j

terminate ? node

else send REJECT(newroot) to j;

// comments -

// if newroot = myroot , then parent is already identified

When ACCEPT(newroot) arrives from j :

if newroot = myroot then

children \leftarrow children $\cup \{j\}$

if (children \cup unrelated) = neighbors \{parent\} then

if i = myroot then terminate

else send ACCEPT(myroot) to parent.

When REJECT(newroot) arrives from j :

if newroot = myroot then

unrelated = unrelated $\cup \{j\}$

if (children \cup unrelated) = (neighbors \{parent\})

then i = myroot then terminate

else

Send ACCEPT(myroot) to parent.

// if newroot < myroot then ignore the message.

// newroot > myroot \rightarrow will never occur -

Asynchronous concurrent-initiator DFS spanning tree algorithm -

Algorithm :-

(local variables)

int parent, myroot $\leftarrow 1$

set of integers children $\leftarrow \emptyset$;

set of integers Neighbors, unknown \leftarrow set of neighbors.

(message types)

QUERY, ACCEPT, REJECT.

When the node wants to initiate the algorithm as a root,

if (parent = 1) then send QUERY(i) to i (itself)

When QUERY(newroot) arrives from j -

if myroot < newroot then

parent $\leftarrow j$;

myroot \leftarrow newroot

unknown \leftarrow set of neighbors.

unknown \leftarrow unknown $\setminus \{j\}$ //excluding j.

if unknown $\neq 0$ then

delete some x from unknown

send QUERY(myroot) to x,

if unknown = 0 then

send ACCEPT(myroot) to j;

else if myroot = newroot then

Send REJECT to j

// if newroot < myroot , ignore the query

when ACCEPT (newroot) or REJECT (newroot) arrives from j :

if newroot = myroot then

if ACCEPT arrived then

children \leftarrow children $\cup \{j\}$

if unknown = 0 then

if parent $\neq i$ then

Send ACCEPT (myroot) to parent.

else if parent = i then

set i as the root
terminate .

delete some x from unknown

Send QUERY (myroot) to x .

// if newroot < myroot , ignore the query since

// sending QUERY to j , i has to update its myroot to j

// will update its myroot to a higher root identifies

// when it a QUERY initiated by it

// newroot > myroot , will never occur .

The Bellman-Ford Shortest path algorithm \Rightarrow

It performs $(n-1)$ times, the relaxation of every edge in the graph.

Input - A weighted directed graph \vec{G} with n vertices and v of \vec{G} .

Output - A label $D[u]$ for each vertex u of \vec{G} such that $D[u]$ is distance from v to u

$$D[v] \leftarrow 0.$$

for each vertex $u \neq v$ of G do

$$D[u] = +\infty.$$

for $i \leftarrow 1$ to $n-1$ do enuring

for each (directed) edge (u, z) from u do

perform the relaxation operation on (u, z)

if $D[u] + w(u, z) < D[z]$ then

$$D[z] \leftarrow D[u] + w(u, z).$$

if there are no edges left with potential relaxation operation

then return the label $D[u]$ of each vertex u

else return \vec{G} contains a negative weight cycle

Single Source Shortest Path algorithm : Synchronous Bellman-Ford

weighted graph : undirectional links

(local variables)

int length $\leftarrow \infty$

int parent $\leftarrow -1$

Set of int Neighbors \leftarrow set of neighbors.

Set of int {weight_{ij}, weight_{ji} | j ∈ Neighbors} \leftarrow the known values of the weights of incident links.

(message types)

UPDATE

if $i = i_0$ then length $\leftarrow 0$

for round $i = 1$ to $n-1$ do

send UPDATE (i , length) to all neighbors.

wait UPDATE (j , length) from each $j \in$ neighbors.

for each $j \in$ Neighbors do

if (length > length_j + weight_{ij}) then

length \leftarrow length_j + weight_{ij} ;

parent = j

→ Source is i_0

→ Code is for process i , $p_i \leq i \leq n$

single source shortest path algorithm: Asynchronous
Bellman Ford

negative cycles \rightarrow process just goes on.

\rightarrow also doesn't allow negative cycles-weight cycles.

(local variables)

int length $\leftarrow \infty$

set of int Neighbors \leftarrow set of neighbors.

set of int {weight_{j,i}, weight_{j,j}} \leftarrow known values of the
weights of incident links.

(message types)

UPDATE

if $i = i_0$ then length $\leftarrow 0$;

send UPDATE ($i_0, 0$) to all neighbors;

terminate;

when UPDATE (i_0, length_j) arrives from j —

if ($\text{length}_j > \text{length}_i + \text{weight}_{j,i}$) then

length $\leftarrow \text{length}_j + \text{weight}_{j,i}$;

parent $\leftarrow j$

send UPDATE (i_0, length) to all neighbors

\rightarrow Depends on the number of hops

algorithm will
definitely
terminate

Distributed File Systems

A distributed file system enables programs to store and access remote files exactly as they do local ones, allowing users to access files from any computer on a network.

[Book -
[oulearis]]

- Performance and reliability experienced for access to files stored at a server should be comparable to that of files stored on local disks.

(Emulates unix file systems at different levels of scaling & performance)

- Sharing of resources → Key goal for distributed systems.
- Sharing of stored information → the most important aspect of distributed resource sharing.
- Design of large scale-wide area read-write file storage systems poses problems of load balancing (few servers might be frequently requested), reliability, availability and security.

Basic distributed file system -

- Objective: Emulate the functionality of a non distributed file system.

Storage systems and their properties →

Sharing, Persistence, Distributed (cache) replicas,
Consistency maintenance

Main Memory

	Sharing	Persistence	Distributed cache replicas	Consistency Maintenance	Example
Main Memory	x	x	x	1	RAM

	Sharing	Persistence	Distributed cache replicas	Consistency Maintenance	Example
File System	x	✓	x	1	UNIX File System

	Sharing	Persistence	Distributed cache replicas	Consistency Maintenance	Example
Distributed File System ↓ (address space is not the same)	✓	✓	✓	✓ (slightly weaker)	SUN NFS

	Sharing	Persistence	Distributed cache replicas	Consistency Maintenance	Example
Web	✓	✓	✓	x each client may modify - it can be modified multiple times	Web Server

	Sharing	Persistence	Distributed cache replicas	Consistency Maintenance	Example
Distributed Shared Memory ↓ address space is the same) - memory is divided into N blocks for N processes	✓	x	✓	✓	

Each processor → one cache. $\xrightarrow{\text{so we need}}$ cache coherent protocols

↓
can be accessed by multiple processes
→ if cache is modified, other processes have to be informed.

<u>Remote objects</u>	✓	✗	✗	1	CORBA
RMI / ORB					

<u>Persistent object store</u>	✓	✓	✗	1	CORBA Persistence State Service
--------------------------------	---	---	---	---	---------------------------------

<u>Peer-to-peer storage system</u>	✓	✓	✓	✗	2 (weaker) (stronger)
------------------------------------	---	---	---	---	-------------------------------------

Types of consistency

1 : strict one copy

✓ : slightly weaker guarantees.

2 : considerably weaker guarantees.

Characteristics of the file system →

- * responsible for the organisation, storage, retrieval, naming, sharing and protection of files.
- They provide a programming interface

File attribute record structure →

- file length
- creation timestamp
- Read timestamp
- write timestamp
- attribute timestamp
- Reference count

} managed by file system.

- Owner
 - file type
 - Access control list
- } Managed by user programs.

file System modules : (for nondistributed file system)

- Directory Module — relates filenames to file IDs.
- File Module — relates file IDs to particular files.
- Access control module — checks permission for operation requested.
- File Access module — reads or writes file data or file attributes
- Block Module — accesses and allocates disk blocks.
- Device Module — performs disk I/O and buffering

Distributed file Systems →

Requirements :-

- Initially, they offered "access transparency" and "location transparency".
 - Later, performance, scalability, concurrency control, fault tolerance and security were added.
- ① Transparency — The design must balance the flexibility and scalability that derive from transparency against software complexity and performance.
1. Access Transparency — client programs should be unaware of the distribution of files.

2. — Location transparency — client programs should see a uniform file namespace
(may be located in diff servers, diff places)
3. — Mobility transparency — Neither client programs nor system administration tables in client nodes need not be changed when the files are moved.
4. — Performance transparency — Client programs should continue to perform satisfactorily while the load on the service varies within a specified range.
5. — Scaling transparency — The service can be expanded by incremental growth to deal with a wide range of nodes and network sizes.

② Concurrent file updates —