

In D.S. processes often compete as well as cooperate to achieve a common goal. It is often required that sites reach mutual agreement.

for ex. In ~~Distribut~~ distributed data base system data managers at sites must agree on whether to commit or to abort a transaction

Reaching an agreement requires that sites have knowledge about the value of other site.

When the system is free from failure an agreement can easily be reached among the processors.

When faulty processes are present in the system; they can send conflicting values to other processors preventing them from reaching an agreement.

In the presence of faults, processors must exchange their value with other processors and delay the values received from other processors several time to isolate the effect of faulty processors

Entire process of reaching an agreement is called an agreement protocol.

In Agreement problem non faulty processor in a D.S. should ~~not~~ be able to reach a common agreement even if certain components in the system are faulty

## System Model

- (1) There are  $n$  processors in the system and at most  $m$  of the processors can be faulty.
- (2) The processors can directly communicate with other processors by msg passing. System is logically fully connected.
- (3) A receiver processor always knows the identity of the sender processor of the message.
- (4) The communication medium is reliable.
- (5) Agreement is to be reached b/w only 2 values 0 and 1.

## Classification of Agreement Problem

(1) Byzantine Agreement Problem

(2) Consensus Problem

(3) Interactive Consistency Problem

## Byzantine Agreement Problem

A single value, which is to be agreed on, is initialized by an arbitrary processor and all non-faulty processors have to agree on that value.

In this problem an arbitrarily chosen processor called the source processor, broadcast its initial value to all other processors.

A soln to the Byzantine agreement problem should meet the following 2 objectives:

Agreement = All non faulty processor agree on the same value ②

Validity = If the source process is non faulty then the common agreed upon value by all non faulty processors should be initial value of the source.

### The Consensus problem

\* Every processor broadcast its initial value to all other processors. Initial value of the processors may be different.

A protocol for reaching consensus should meet the following condition.

Agreement: All non faulty processor agree on the same single value

Validity: If the initial value of every non faulty processor is  $v$ , then agreed upon common value by all non faulty processor must be  $v$

### The interaction consistency Problem

Every processor broadcast its initial value to all other processors. The initial value of the processor may be different.

A protocol for the interaction consistency problem should meet the following condition.

Agreement: All non faulty processor agree on the same vector  $(v_1, v_2, \dots, v_n)$

Validity: If the  $i^{th}$  processor is non faulty and its initial value is  $v_i$  then the  $i^{th}$  value to agreed on by all non faulty processor must be  $v_i$

## Solution for Byzantine Agreement Problem

- ⇒ First defined and solved by Lamport
- ⇒ Some broadcast its initial value to all other processors.
- ⇒ Processors send their value to other processors and also relay received values to other
- ⇒ During execution faulty processors may confuse by conflicting values
- ⇒ If faulty processors dominate in no. they can prevent non faulty processors from reaching an agreement
- ⇒ No of faulty processors should not exceed certain limit.
- ⇒ Pease showed that in a fully connected network it is impossible to reach an agreement if no faulty processor 'm' exceeds  $(n-1)/3$ .  
 $n = \text{no of processors}$

## Lamport Shostak Pease Algorithm:

This algo also known as oral message algorithm OM( $m$ ) where  $m$  is the no of faulty processor.  
 $n = \text{no of processors}$  and  $n \geq 3m+1$

Algo is recursively defined as follow.

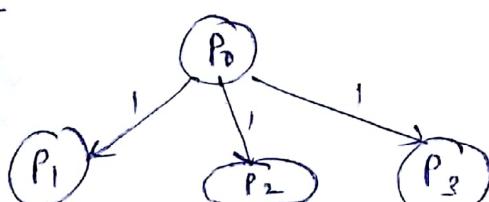
Alg. OM(0)

Source processor send its value to every processor  
Each processor ~~sends~~ uses the value it receives from source. If no value is received default value 0 is used

(3)

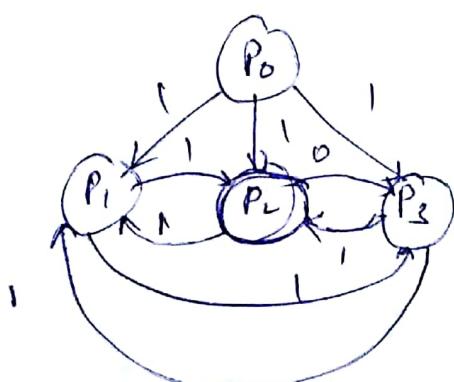
Alg OOM(m) m&gt;0

- 1 The source processor sends its value to every processor.
- 2 for each  $i$ , let  $v_i$  be the value processor  $i$  receives from source.
- 3) Processor T act as the new source and initiates algo OM(m-1) when it sends the value  $v_i$  to each of the  $(n-2)$  other processors.
- 4) for each  $i$  and  $j$  let  $v_{ij}$  be the value processor  $i$  received from processor  $j$  in Step 3  
 Processor T use the value majority  $(v_1, v_2, \dots, v_{n-1})$   
 The function majority  $(v_1, v_2, \dots, v_{n-1})$  computes the majority value if exists otherwise it uses default value 0

Example

$\Rightarrow P_2 = \text{faulty process}$   
 $P_0 = \text{non faulty.}$

Process  $P_0$  executes the algo OM(1)



Process  $P_1, P_2, P_3$  execute the algo OT1(0)

$$P_1 \text{ received} = (111)$$

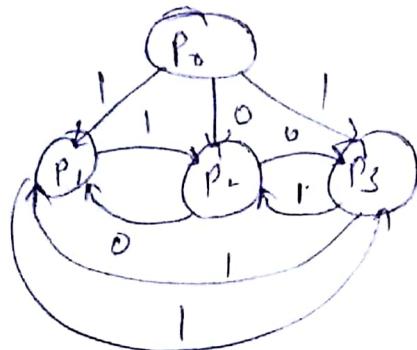
$$P_2 = (11\cancel{1})2$$

$$P_3 = (1\cancel{1}0)$$

$$\text{Majority value} = 1$$

both Cond<sup>n</sup> of Byzantine Agreement are satisfied

In  
m  
Ph  
Se



$P_0 = \text{faulty}$   
Process  $P_0, P_1, P_2, P_3$  execute OM(1)  
& OM(0)

$$P_1 \text{ received} = \{1, 0, 1\}$$

$$P_2 = \{0, 1, 1\}$$

$$P_3 = \{1, 0, 1\}$$

$$\text{Majority} = 1$$

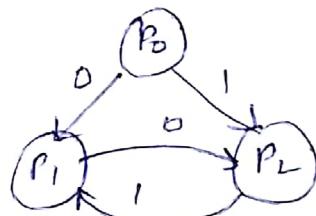
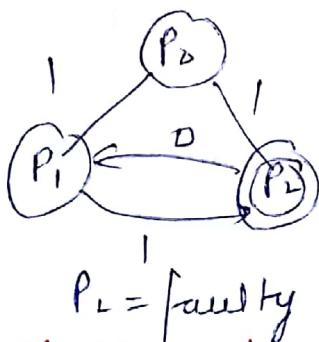
Non-faulty process or agrees on same value.

Byzantine Agreement can not be reached among three processes if 1 process is faulty

$P_0$  is some processor

$P_0$  is Non faulty

$P_0$  is faulty



### Application of Agreement Algo.

1) Fault Tolerance Clock Synchronization

D.S. requires physical clocks to synchronize  
Physical clocks have drift problem

Agreement protocol may help to reach a common clock value

2) Atomic Commit in DBS

DBS sites must agree whether to commit & abort the transaction

Agreement protocol may help to reach a consensus

In the problem of atomic Commit. sites of a DBS must agree whether to commit or abort a transaction.

Phase First = sites execute their part of a distributed transaction and broadcast their decisions to all other sites

Second Phase = Each site based on what it received from other sites in the 1st phase decides whether to commit or abort its part of the distributed transaction

⇒ Every site receives an identical response from all other sites. They will reach the same decision. If some sites behave maliciously they can send a conflicting response to other sites causing them to make conflicting decisions. In this situation we can use algorithm for the Byzantine Agreement to insure that all non-faulty processor reach a common decision about distributed transaction.

It works as follow:

In the 1st phase after a site has made a decision it starts the Byzantine Agreement

In the 2nd Phase = Processors determine a common decision based on the agreed vector of Phase

# Distributed File Systems.

①

**Introduction:** A distributed file system is a resource management component of a distributed operating system.

It implements a common file system that can be shared by all the autonomous computers in the system.

Two important goals of D.F.S.

→ **M/w Transparency:** User don't have to be aware of the location of files to access them. This property of a D.F.S. is known as M/w transparency.

→ **High Availability:** Another major goal of D.F.S is to provide high availability.

System failure or regularly scheduled activities such as back ups or maintenance should not result in the unavailability of files.

## Architecture

- \* In D.F.S. files can be stored at any machine or the computation can be performed at any machine.
- \* When a machine needs to access a file stored on remote machine, the remote machine performs the necessary file access operations and return data if a read operation is performed.

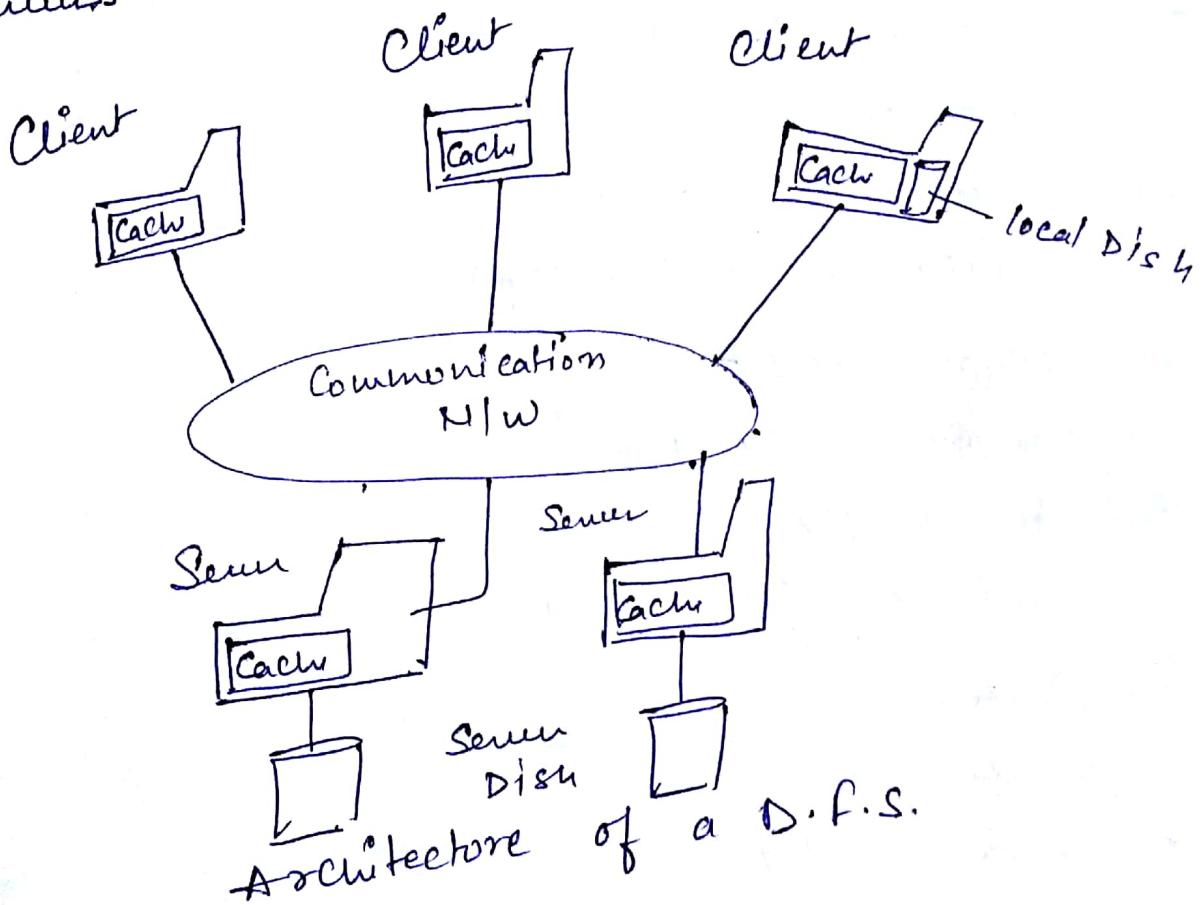
The two most important services present in a D.F.S. are:

**Name Service:** is a process that maps names specified by client to stored objects such as files and directories.

Cache Manager: is a process that implements file caching. (2)

In file caching, a copy of data stored at a remote file server is brought to the client's machine.

Cache manager can be present at both client and file servers



# Mechanisms for Building D.F.S.

(3)

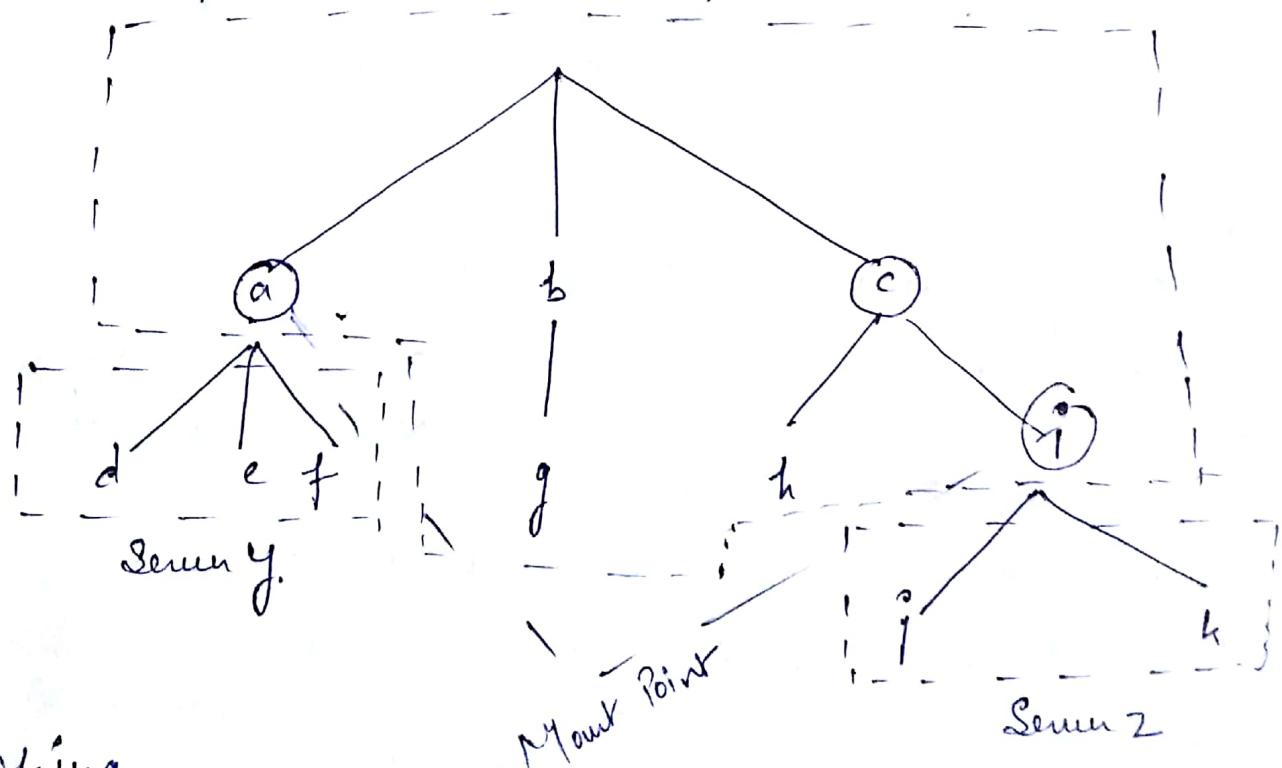
## Mounting:

- Mounting : Binding together different file name spaces to form a single hierarchically structured name space
- A name space can be bounded to or mounted at external node or a leaf node of a name space tree.

## Mount Table

- \* The kernel maintains a structure called the mount table which maps mount point to appropriate storage device.

Sunux



- ## Caching
- = A common technique to exploit locality and reduce delay.
  - In file caching, a copy of data stored at remote file server is brought to the client when referenced by the client.

Hints = Caching without Cache Coherence Cache Coherence operations too expensive in D. Systems.

Hints work well for applications that can recover from individual data.

Bulk Data Transfer = Transferring data in bulk reduces the protocol processing overhead at both server and clients.

Encryption: Used to enforcing security in D.S.

## Design issues

### (1) Naming and Name Resolution

\* A Name in file system is associated with an object  
Name resolution refers to the process of mapping a name to an object

Name Space - a collection of names

Three different methods to name file.

- (1) Concatenate the host name to the name of files that are stored on that host.
- filename is unique
  - No NW transparency
  - Moving a file from machine to another may require modifications to the application
- Advantage - Name resolution is simple. - Don't need a name server

## (2) Mount Remote directories onto local Directory

- Once mounted, accessing the files becomes location transparent  
Name resolution is simple

## (3) Have a single Global directory: all files in the system belongs to a single Name space

- Not scalable, how to maintain a systemwide unique file names.

NameServer: In a centralized system, name resolution can be accomplished by maintaining a table that maps names to objects.

In distributed systems, name servers are responsible for name resolution

A Name server is a process that maps names to stored object

\* The easiest approach to name resolution in D.S. is for all client to send their queries to a single name server which maps name to objects.

→ Drawback: If the name server crashes the entire system is drastically affected.

→ (a) Degrade the performance of the system

\* Second Approach → Having several name servers wherein each server is responsible for mapping object stored in different domain

## Caches on Disk or Main Memory:

Advantage for caching on main memory -

- Can be used in disk less workstation
- accessing a cache in memory is much faster
- Same technique for server cache and client cache

Disadvantage - Compete with virtual memory

- Virtual memory system more complex

Data cannot be in both VM & in cache

large file cannot be cached completely.

Advantage of caching on Disk:

large file can be cached. Virtual memory is simple

Writing Policy : It decides when a modified cache block at a client should be transferred to the server

The simplest Policy - Write Through - In write through all writes requested by the applications at clients are also carried out at the server immediately.

The main advantage of write through = reliability  
disadvantage = In the event of client crash, little info. is lost. Does not take advantage of cache

Delayed Writing Policy - Delays the writing at the server

Modifications due to a write are reflected at the server after some delay.

Access Locality: Some data can be deleted right away - 20-30%. Data are erased within 30 seconds.

May loss data when Client Crashes.

⇒ Another writing Policy delays the updating (s) of the files at the server until the file is closed at the client.

The traffic at the server depends on the avg. period that files are open.

If the avg period for which files are open is short, then this policy does not greatly benefit from delaying the updates.

If the avg period for which files are open is long this policy is also susceptible to losing data in the event of a client crash.

## Cache Consistency

There are 2 approaches to guarantee that the data returned to the client is valid.

(1) Server-initiated: Server informs Cache managers whenever the data in the client cache become stale. Cache Manager at client can then receives the new data or invalidate the blocks containing the old data in their cache.

(2) Client-initiated - the client make sure the cache is clean before returning the data to the application.

\* Both the schemes are expensive

\* Not allow caching when concurrent write sharing occurs.

\* Sequential - writes Sharing may also cause problem a client opens a file that has recently modified and closed by another client.

- ⇒ Out date blocks in the cache: timestamp the cache copy. Contacting the server can find the inconsistency.
- ⇒ Update data in another client's cache not in the server: flushing writes when other clients open

## Availability

- (\*) Replication is the main mechanism
- (\*) How to keep replicas consistent
- (\*) How to detect inconsistencies among replicas
- (\*) Consistency problem may decrease the availability
- (\*) Replica Management: Voting mechanism to read and write to replica

- ## Scalability
- \* How to meet the demand of a growing system?
  - \* Biggest hurdle: Consistency issue.
    - Caching: Reduce H/W load & Server load.
    - Server initiated cache invalidation to maintain cache consistency complexity increases as the system size increases.

(6)

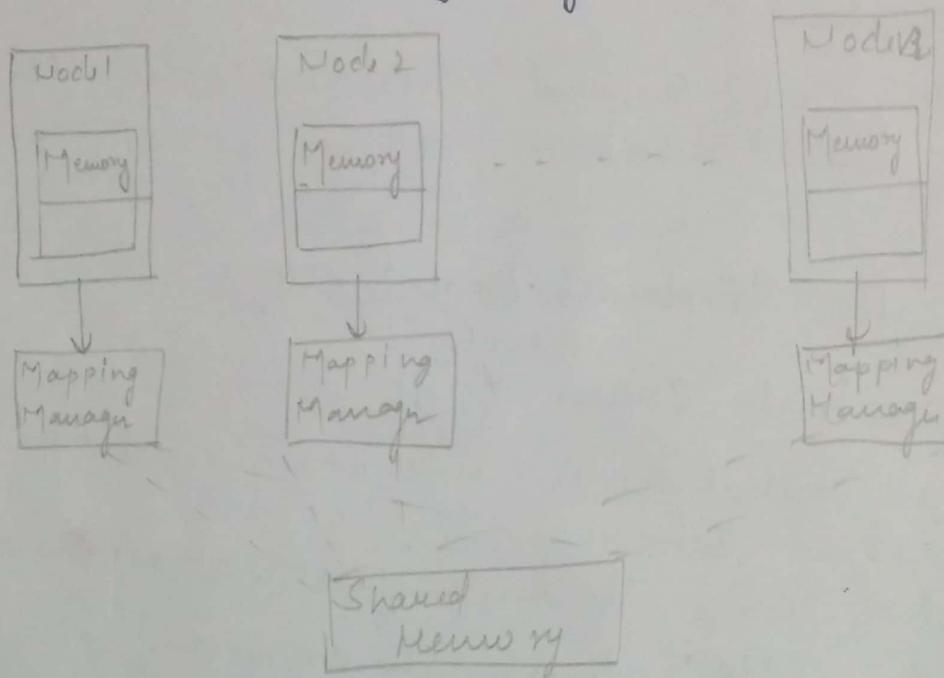
## Semantics

what a user want? Strict Consistency  
users can usually tolerate a certain  
degree of errors in file handling - no  
need to enforce strict consistency

# The Read Replication Algo;

## Distributed Shared Memory.

DSM is a resource management component of a distributed operating system, that implements the shared memory model in distributed system which have no physically shared memory.



## Advantages of DSM.

- Programming is a lot easier
  - No need to deal with communication details.
  - Easy to handle complex data structures
- DSM systems are much ~~expensive~~ cheaper than highly coupled multiprocessor system
- DSM takes advantages of the memory reference locality. - Data are moved in the unit of pages
- DSM can form a large physical memory

- Programs for shared Memory multiprocessors can easily be ported to DSM.

## Challenges in DSM.

- ⇒ How to keep track of the location of remote data?
- ⇒ How to overcome the communication delay & high overhead?
- ⇒ How to allow controlled concurrent access to shared data?

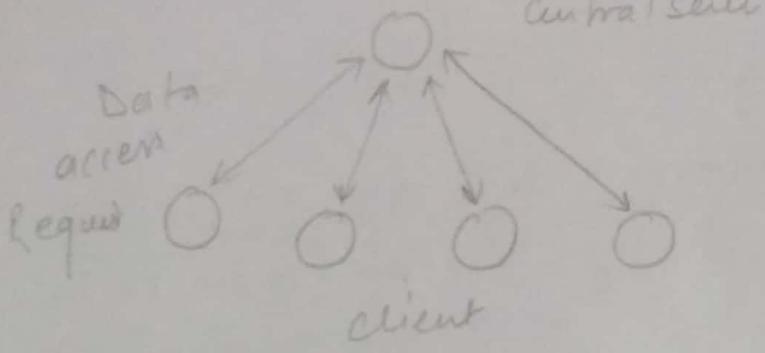
## Algorithm for Implementing DSM.

### (1) The central-Server Algo

- (•) A central server maintains all the shared data.
- (•) for Read: The server just return the data
- (•) for write: Update the data and send acknowledgement for distributed applications.

• A simple working solution to provide shared memory for distributed application.

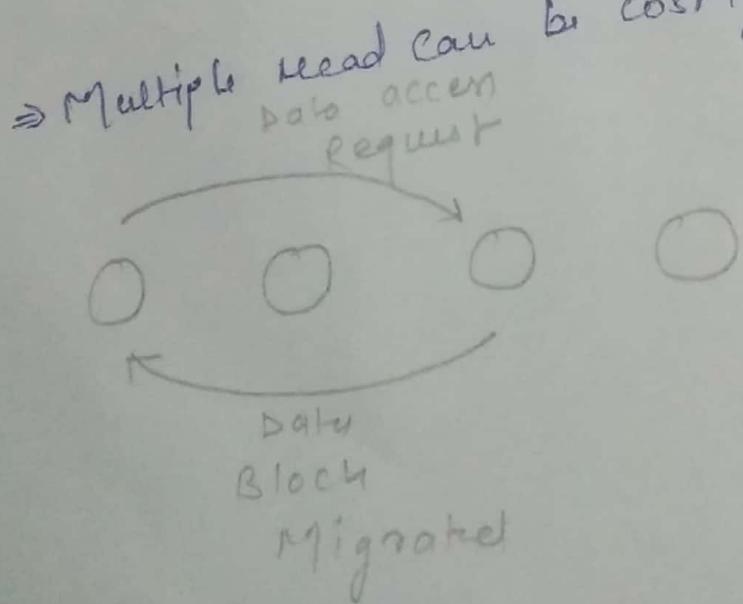
- Low efficiency. -Bottleneck at the server, long memory access latency
- Data can be distributed. - need a directory to store the location of a page



## The migration Algorithms:

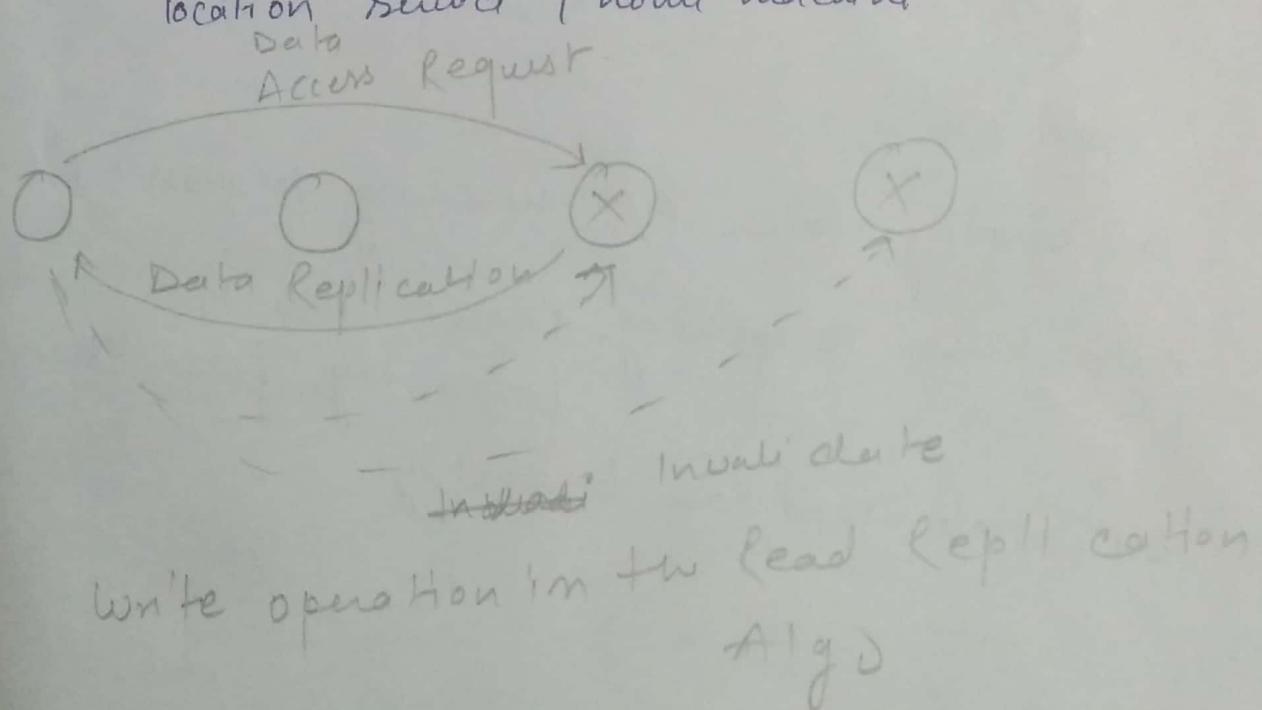
- ⇒ Data is shipped to the location of the data access request. - subsequent accesses are local.
- ⇒ The migration algo allows only one node to access a shared data at time.
- ⇒ For both read/write: get the remote page to the local machine, then perform the operation.

- ⇒ keeping track of memory location: Location service
- ⇒ Problems: thrashing - Pages move b/w nodes frequently. false sharing can be costly



## The Read Replication Algo:

- ⇒ Extends the migration algo by replicating data blocks and allowing multiple nodes to have read access or own node to have read+write access.
- ⇒ On Read, a page is replicated. Mark the page as multiple reader
- ⇒ On write, all copies except one must be updated or invalidated
- ⇒ Multiple read on write
- ⇒ Allowing multiple readers to a page
- ⇒ All the locations must be kept track of:  
location service / home machine



## The Full Replication Algo

- ⇒ Allow multiple nodes to have both read and write access to shared data blocks  
(Multiple reader multiple writer protocol)
- ⇒ Many node can write shared data concurrently  
the access to shared data must be controlled  
to maintain its consistency.
- ⇒ All nodes wishing to modify shared data  
will send the modification to a sequencer
- ⇒ The sequencer will assign a sequence no  
and multicast the modification with the  
sequence no to all the nodes that have  
a copy of the shared data item  
each node processes the modification request  
in the sequence no. order

