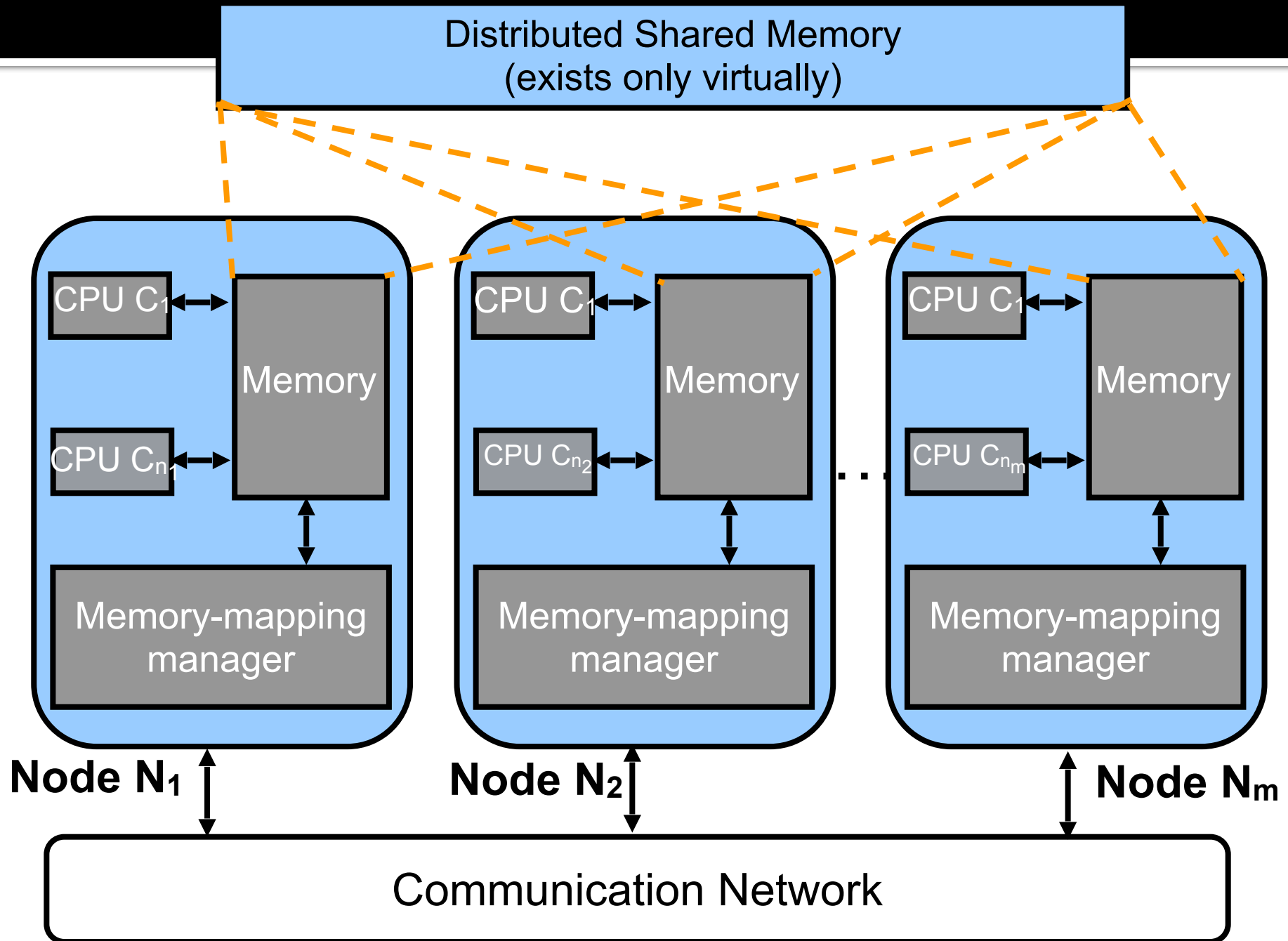


Advanced Operating Systems: Distributed Shared Memory

Motivation

- RPC allows us to pass messages to the processes in the distributed systems.
- RMI allows us to call procedures in the distributed systems.
- We used to have shared memory in uni-processor systems to share data between process.
- It is popular to use shared-memory in tightly-coupled multi-processor systems.
- How about loosely coupled distributed systems?

Distributed Shared Memory (DSM)



Advantages of DSM

- **Simpler Abstraction**
 - Programming distributed memory machines
 - Message passing models is tedious and error prone.
 - Under RPC and message passing, it is difficult to pass context-related data or complex data structures.
- **Better Portability of Distributed Application Programs**
 - Consistent access protocol makes it easier to transit from sequential applications to distributed applications.
 - Migrating shared-memory multiprocessor applications to distributed systems with distributed shared memory is seamless.
- **Better Performance of Some Applications**
 - Locality of Data
 - On-demand data movement
 - Larger memory space
- **Flexible Communication Environment**
- **Ease of Process Migration**

Design and Implementation Issues of DSM

- Granularity: block vs. page
- Structure of shared-memory space
- Memory coherence and access synchronization (consistence)
- Data location and access
- Replacement strategy
- Thrashing
- Heterogeneity

Coherence vs. Consistency

- Coherence concerns only *one* memory location
- Consistency concerns for *all* locations
- A memory system is coherence if
 - it can serialize all operations to that location
 - operations performed by any core appear in program order.
 - it reads return values written by last store to that location.
- A memory system is consistent if
 - if follows the rules of its memory model
 - operations on memory location appears in some defined order.

Coherence vs. Consistency

- Name a few coherence protocol:
 - **Snooping**: snooping is a process where the **individual caches monitor** address lines for accesses to memory locations that they have cached. When a write operation is observed, the cache controller **invalidates** its own copy of the snooped memory location.
 - **Snarfing**: a cache controller watches both address and data in an attempt to update its own copy of a **memory location** when a second master modifies a location in main memory. When a write operation is observed to a location that a cache has a copy of, the cache controller **updates** its own copy of the snarfed memory location with the new data.
- Name a few consistency protocol:
 - **Strict consistency**: if a process reads any memory location, the value returned by the read operation is the value written by the most recent write operation to that location.
 - **Sequential consistency**
 - **Processor consistency**

Coherent but not consistent

- Can you find a memory trace which is coherent but not consistent?

initially A=B=0

process 1

store A := 1

store B := 1

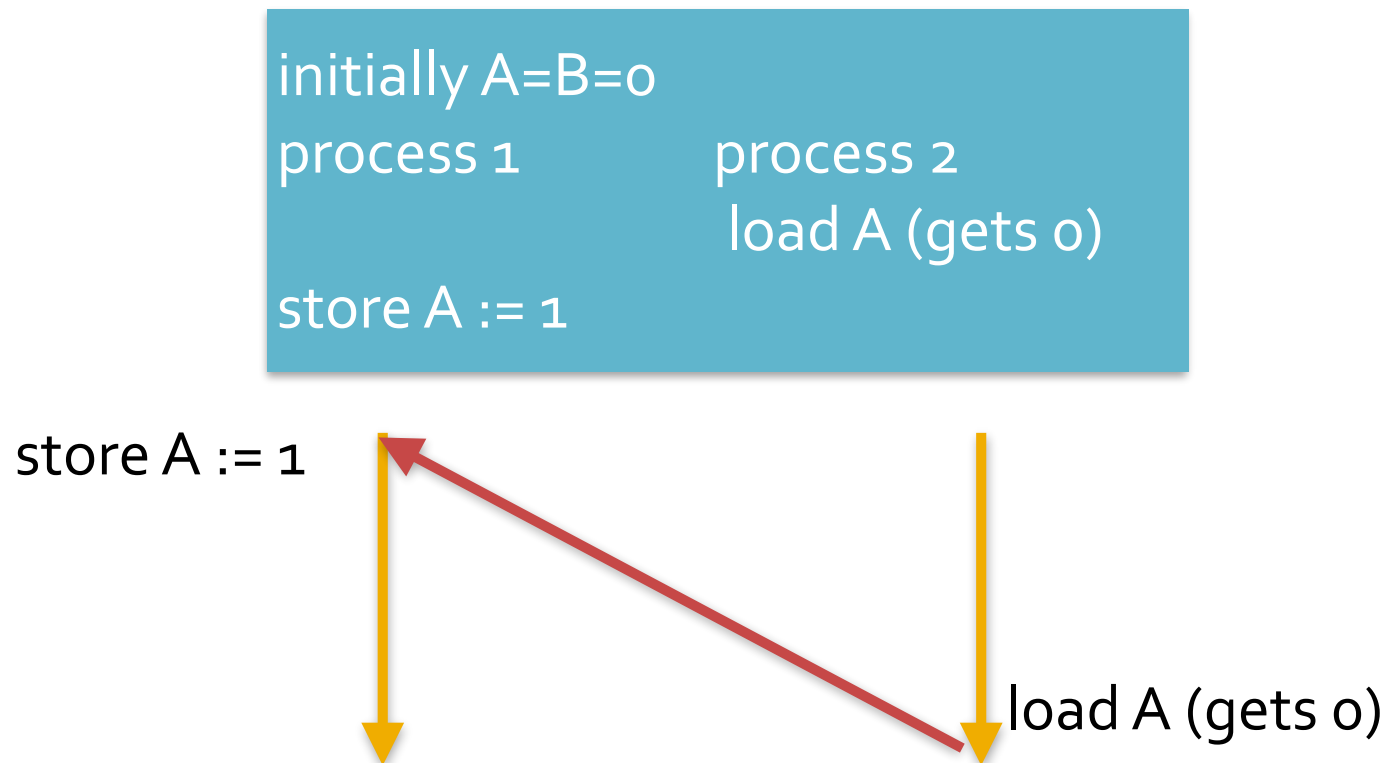
process 2

load B (gets 1)

load A (gets 0)

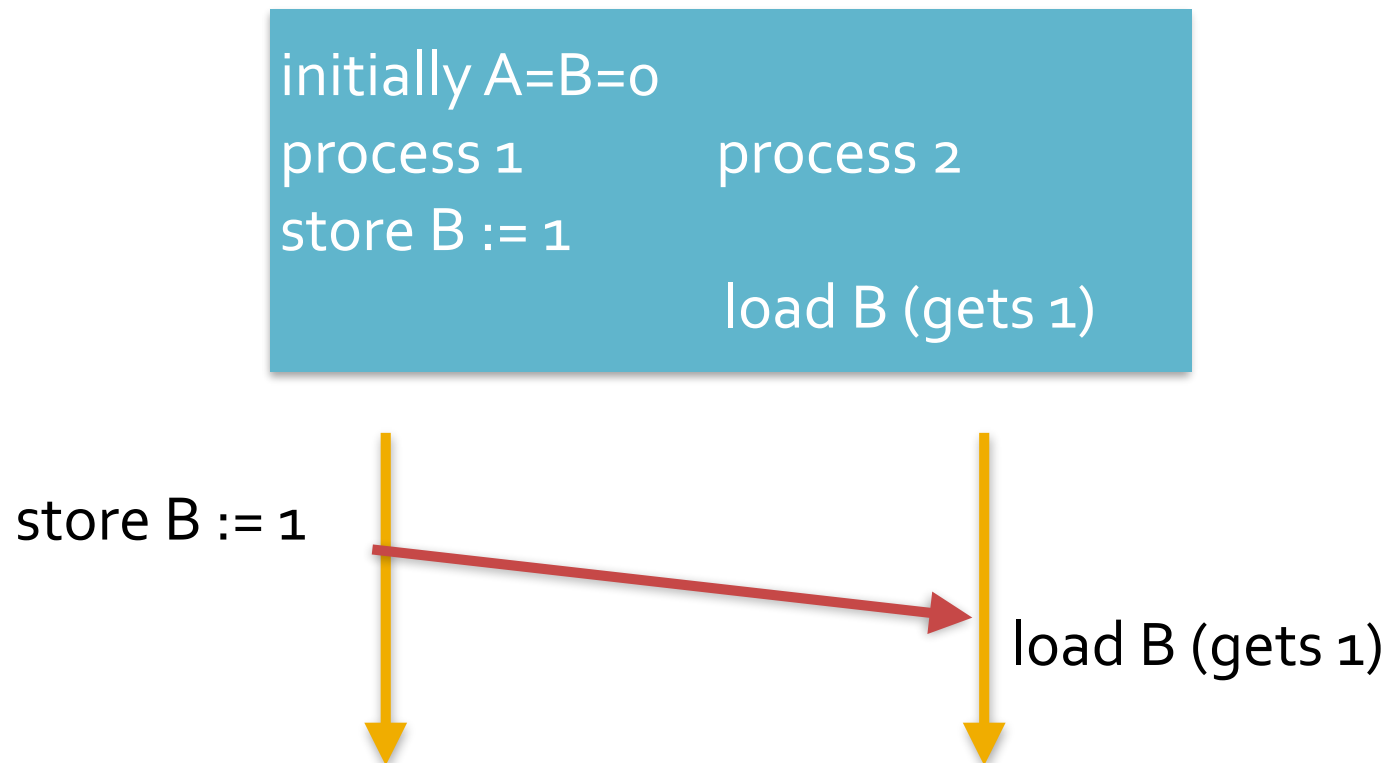
Coherent but not consistent

- Can you find a memory trace which is coherent but not consistent?



Coherent but not consistent

- Can you find a memory trace which is coherent but not consistent?



Coherent but not consistent

- Can you find a memory trace which is coherent but not consistent?

initially A=B=0

process 1

store A := 1

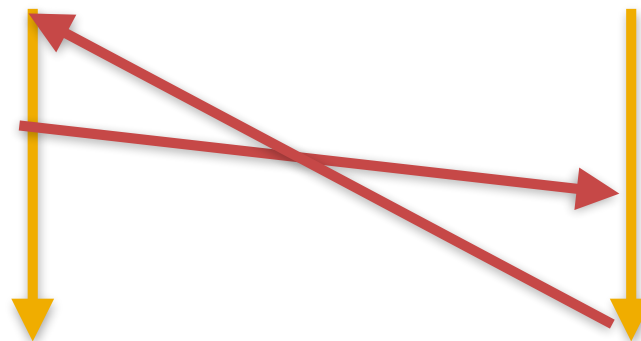
store B := 1

process 2

load B (gets 1)

load A (gets 0)

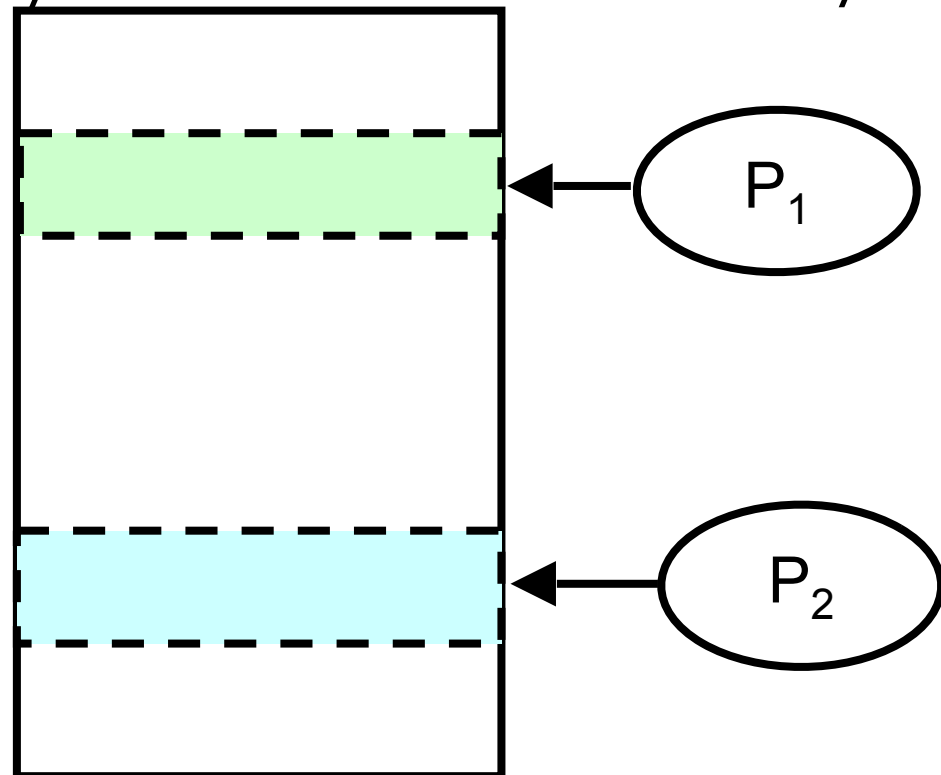
store A := 1
store B := 1



load B (gets 1)
load A (gets 0)

Granularity – How to select block size

- Block: the unit for transmitting data.
- Trade-off: network traffic vs. parallelism
- What's the difference between multi-processor system and distributed systems in terms of memory access?
- Factors to consider:
 - Paging overhead
 - Directory size
 - Thrashing
 - False sharing



One data block

Using page size as block size

- The system can use existing page fault schemes.
- The system can use existing access right control.
- If a page can be fitted into a packet, page sizes do not impose undue communication overhead.
- A page size is a suitable data entity with respect to memory contention.

Structure of Shared-Memory Space

- Structure: the abstract view of the shared-memory space
 - One may see the DSM as a storage of words and
 - The other may see the DSM as a storage of data objects.
- It is related to the choice of block size.
- Three common structures:
 - No structuring
 - Fixed grain size for all applications
 - Easier to choose any suitable page size as the unit of sharing
 - Structuring by data type
 - Variable grain size
 - Complicated design and implementation
 - Structuring as a database
 - Tuple space: memory ordered by their content.
 - Accessed by specifying the number of their fields and their values via special access functions
- How does the type of structure affect the implementation of your systems?

Consistency Models

- Consistency models: the degree of consistency that has to be maintained
- Ongoing researches: relax the requirements to a greater degree.
- Example of different consistency models:

看股票

看開票

- Which one aims on **ordering**?
- Which one aims on **results**?

Consistency Models

- Stronger consistency model vs. weaker consistency model
- Available models:
 - Strict consistency model
 - Sequential consistency model
 - Causal consistency model
 - Pipelined random-access memory consistency model
 - Processor consistency model
 - Weak consistency model
 - Release consistency model

Strict Consistency Model

- The value returned by a read operation on a memory address is always the same as the value written by the most recent write operation to that address.
- All writes instantaneously become visible to all processes.
- What you need:
 - read/write operations must be correctly ordered
 - an absolute global clock

Consistency Models – Strict Consistency

Node N_1

Node N_2

Node N_3

```
{
  ...
  a=d
  c=a+b
  a=4
  ...
}
```

```
{
  ...
  a=10
  ...
  print(a)
}
```

```
{
  ...
  e=foo()
  f=bar()
  ...
}
```

Node N_1

Node N_2

Node N_3

Strict Consistency Model

Global clock

...

$r_1(d)$

...

$w_1(a)$

$r_1(a)$

$r_1(b)$

...

$w_1(c)$

$w_1(a)$

...

...

$w_2(a)$

...

$r_2(a)$

...

$w_3(e)$

...

$w_3(f)$

...

...

$r_1(d)$

$w_2(a)$

$w_1(a)$

$r_1(a)$

$r_1(b)$

$w_3(e)$

$w_1(c)$

$w_1(a)$

$w_3(f)$

$r_2(a)$

...

...

$r_1(d)$

$w_2(a)$

$w_1(a)$

$r_1(a)$

$r_1(b)$

$w_3(e)$

$w_1(c)$

$w_1(a)$

$w_3(f)$

$r_2(a)$

...

...

$r_1(d)$

$w_2(a)$

$w_1(a)$

$r_1(a)$

$r_1(b)$

$w_3(e)$

$w_1(c)$

$w_1(a)$

$w_3(f)$

$r_2(a)$

...



Sequential Consistency Model

- It was proposed by [Lamport](#) in '79.
- All processes see the same order of all memory access operations on the shared memory.
 - The orders seen by processes must be the same but
 - They are not necessary to be equal to the **EXACT** orders.
- The sequential consistency model does not guarantee that a read operation on a particular memory address always returns the same value as written by the most recent write operation to that address.
- Running a program twice may not give the same result in the absence of explicit synchronization operations.
- A sequential consistent memory provides *one-copy/single-copy* semantics because all the processes sharing a memory location always see the same contents.



MORE ACM AWARDS

A.M. TURING AWARD



Search

TYPE HERE



A.M. TURING AWARD WINNERS BY...

ALPHABETICAL LISTING

YEAR OF THE AWARD

RESEARCH SUBJECT



PHOTOGRAPHS

BIRTH:

7 February 1941 in New York, New York

EDUCATION:

Bronx High School of Science (1957); B.S. (Massachusetts Institute of Technology, Mathematics, 1960); M.S. (Brandeis University, Mathematics, 1963); PhD (Brandeis University, Mathematics, 1972).

EXPERIENCE:

Massachusetts Computer Associates, 1970-1977; SRI International, 1977-1985; Digital Equipment Corporation and Compaq, 1985-2001; Microsoft Research, from 2001.

HONORS AND AWARDS:

Dijkstra Prize for the paper "Time

LESLIE LAMPORT

United States – 2013

CITATION

For fundamental contributions to the theory and practice of distributed and concurrent systems, notably the invention of concepts such as causality and logical clocks, safety and liveness, replicated state machines, and sequential consistency.



SHORT ANNOTATED
BIBLIOGRAPHY



ACM DL
AUTHOR PROFILE



ACM TURING AWARD
LECTURE VIDEO



RESEARCH
SUBJECTS

*If we could travel back in time to 1974, perhaps we would have found Leslie Lamport at his busy local neighborhood bakery, grappling with the following issue. The bakery had several cashiers, but if more than one person approached a single cashier at the same time, that cashier would try to talk to all of them at once and become confused. Lamport realized that there needed to be some way to guarantee that people approached cashiers one at a time. This problem reminded Lamport of an issue which has been posed in an earlier article by computer scientist Edsger Dijkstra on another mundane issue: how to share dinner utensils around a dining table. One of the coordination challenges was to guarantee that each utensil was used by at most one diner at a time, which came to be generalized as the **mutual exclusion** problem, exactly the challenge Lamport faced at the bakery.*

One morning in 1974, an idea came to Lamport on how the bakery customers could solve mutual exclusion among themselves, without relying on the bakery for help. It worked roughly like this: people choose numbers when they enter the bakery, and then get served at the cashier according to their number ordering. To choose a number, a customer asks for the number of everyone around her and chooses a number higher than all the others.

This simple idea became an elegant algorithm for solving the mutual exclusion problem without requiring any lower-level indivisible operations. It also was a rich source of future ideas, since many issues had to

Consistency Models – Sequential Consistency

Node N ₁	Node N ₂	Node N ₃
{ ... ä=d c=a+b a=4 ... }	{ ... ä=10 ... print(a) }	{ ... ë=foo() f=bar() ... }
...		
r ₁ (d)		
...	w ₂ (a)	...
w ₁ (a)	...	
r ₁ (a)		
r ₁ (b)		
...		w ₃ (e)
w ₁ (c)		
w ₁ (a)		...
...	r ₂ (a)	w ₃ (f)
		...

Sequential Consistency Model

Global clock

Node N ₁	Node N ₂	Node N ₃
...
r ₁ (d)	r ₁ (d)	r ₁ (d)
w ₁ (a)	w ₁ (a)	w ₁ (a)
w ₂ (a)	w ₂ (a)	w ₂ (a)
r ₁ (a)	r ₁ (a)	r ₁ (a)
r ₁ (b)	r ₁ (b)	r ₁ (b)
w ₃ (e)	w ₃ (e)	w ₃ (e)
w ₁ (c)	w ₁ (c)	w ₁ (c)
w ₁ (a)	w ₁ (a)	w ₁ (a)
w ₃ (f)	w ₃ (f)	w ₃ (f)
r ₂ (a)	r ₂ (a)	r ₂ (a)
...



What're the difficulties of implementing consistency model?

- Each node/process needs to know which instructions are issued by other nodes/processes.
 - Communications or synchronizations are required among the nodes/processes.
 - Communications/synchronizations will slow down or block the progress.
- Consequently, the performance of the systems become poor.
 - When the number of nodes/processes increase, the penalty increases (exponentially).

Further relaxing the model to avoid communication overhead

- The outcome of a sequence of memory operations depend on
 - execution order and
 - what else?

```
a=0;b=0;  
a=1;  
b=a+2;  
print("a: %d, b:%d\n", a, b);
```

```
a=0;b=0;  
b=a+2;  
a=1;  
print("a: %d, b:%d\n", a, b);
```

Causality 因果關係

Further relax the model

- The outcome of a sequence of memory operations are related to
 - execution order and
 - what else?

```
a=0;b=0;  
a=1;  
b=2;  
print("a: %d, b:%d\n", a, b);
```

```
a=0;b=0;  
b=2;  
a=1;  
print("a: %d, b:%d\n", a, b);
```

When there is no causality, the execution order has no effects.

Causally Related

A memory reference operation (read/write) is said to be potentially causally related to another memory reference operation if the second one might have been influenced in any way by the first one.

```
foo () {
```

```
...
```

```
read(a) ;
```

```
b = a * c;
```

```
write(b) ;
```

```
...
```

```
}
```

Causally related

```
foo () {
```

```
...
```

```
read(a) ;
```

```
b = a * c;
```

```
write(b) ;
```

```
...
```

```
}
```

```
bar () {
```

```
...
```

```
read(d) ;
```

```
e = d * c;
```

```
write(e) ;
```

```
...
```

```
}
```

Not causally related

Casual Consistency Model

- It is proposed by Hutto and Ahamad in '90.
- In the casual consistency model,
 - all processes see only those memory reference operations in the same order that are potentially causally related,
 - memory reference operations that are not causally related may be seen by different processes in different orders.
- A shared memory system is said to support the causal consistency model if all write operations that are potentially causally related are seen by all processes in the same (**correct**) order.
 - Suppose W_2 is causally related to W_1 , i.e., W_2 depends on the results of W_1 .
 - Only (W_1, W_2) is correct. (W_2, W_1) is not.

Consistency Models – Casual Consistency

Node N ₁	Node N ₂	Node N ₃
{	{	{
... a=d c=a+b a=4 ... }	... a=10 print(a) }	... e=foo() f=bar() ... }

...		
r ₁ (d)
...	w ₂ (a)	
w ₁ (a)	...	
r ₁ (a)		
r ₁ (b)		w ₃ (e)
...		
w ₁ (c)		
w ₁ (a)		w ₃ (f)
...	r ₂ (a)	...

Node N ₁	Node N ₂	Node N ₃
---------------------	---------------------	---------------------

Casual Consistency Model

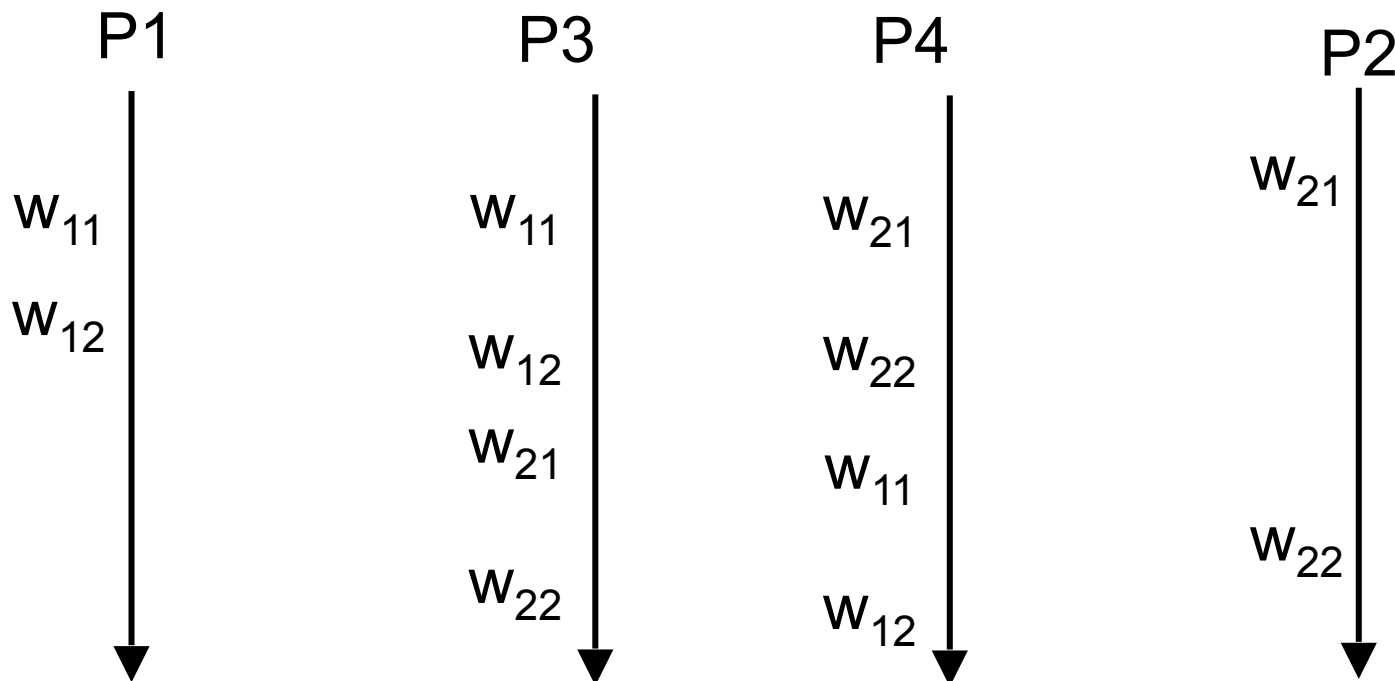
Global clock

...
r ₁ (d)	r ₁ (d)	r ₁ (d)
w ₃ (e)	w ₂ (a)	w ₁ (a)
w ₂ (a)	w ₃ (e)	w ₂ (a)
w ₁ (a)	w ₁ (a)	r ₁ (a)
r ₁ (a)	r ₁ (a)	r ₁ (b)
r ₁ (b)	r ₁ (b)	w ₃ (e)
w ₁ (c)	w ₁ (c)	w ₃ (f)
w ₃ (f)	w ₁ (a)	w ₁ (c)
w ₁ (a)	w ₃ (f)	w ₁ (a)
r ₂ (a)	r ₂ (a)	r ₂ (a)
...



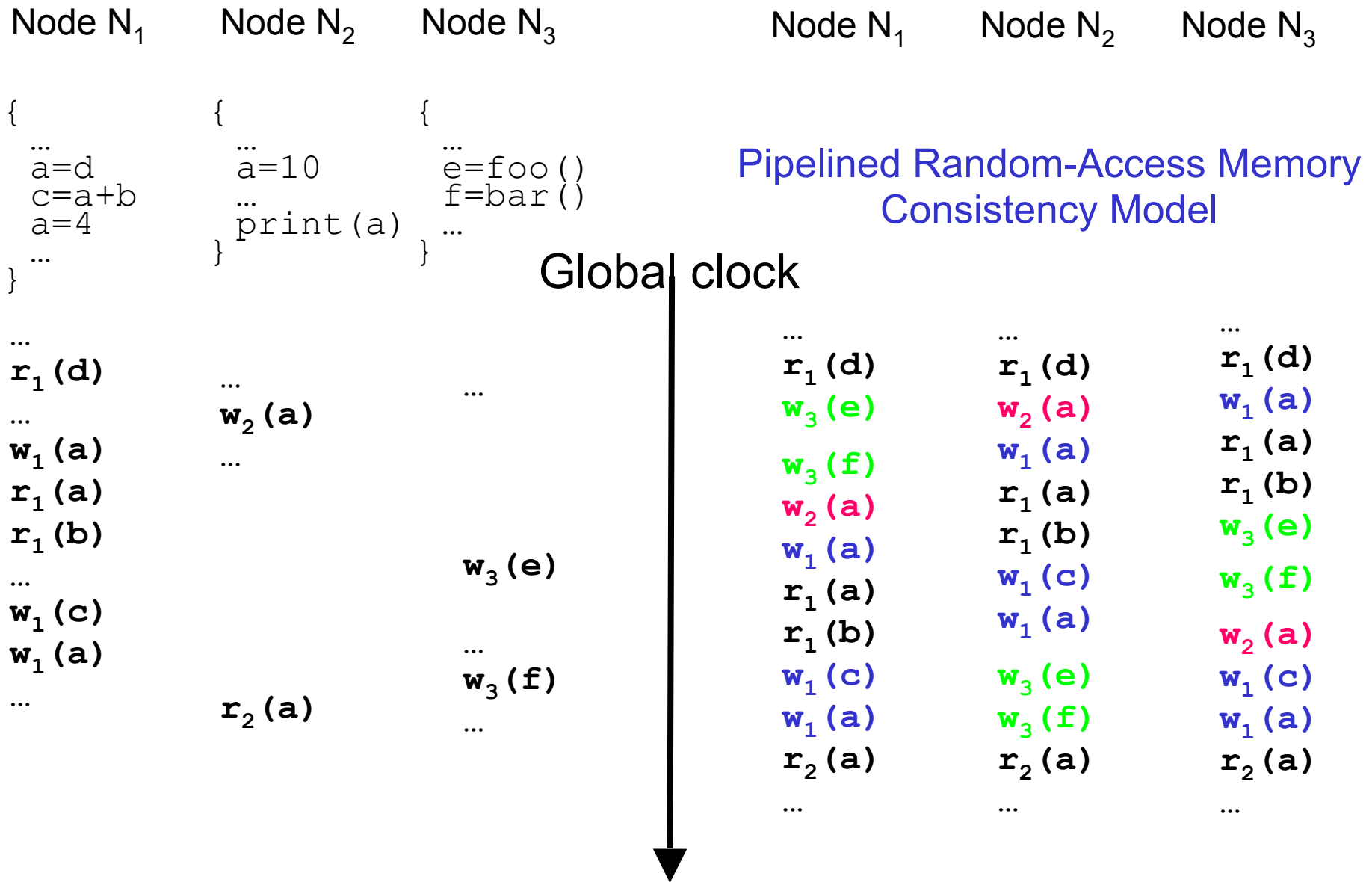
Pipelined Random-Access Consistency Model

- It is proposed by Lipton and Sandberg in '88.
- PRAM Consistency Model:
 - All write operations performed by a single process are seen by all other processes in the order in which they were performed as if all the write operations performed by a single process in a pipeline.
 - Write operations performed by different processes may be seen by different processes in different orders.



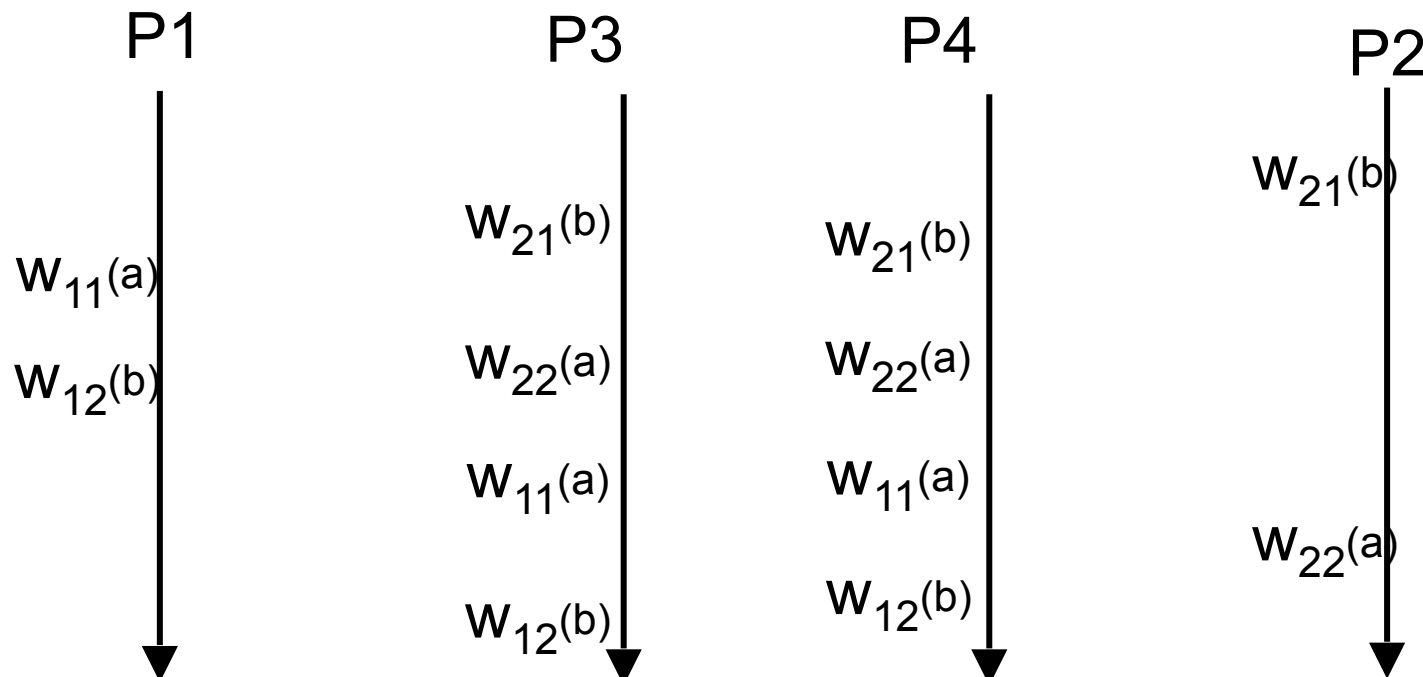
- PRAM Consistency Model is simple and easy to implement.

Consistency Models - PRAM



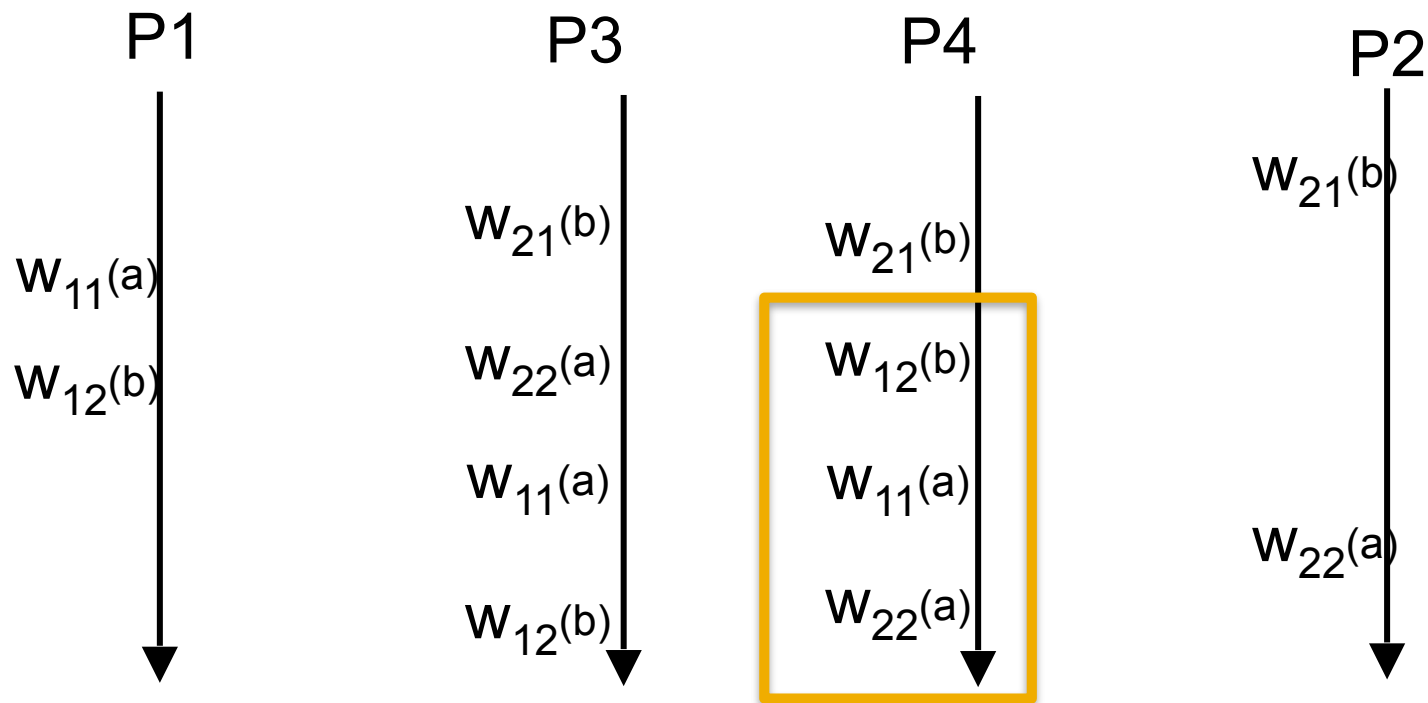
Consistency Models – Processor Consistency Model

- Proposed by Goodman in '89.
- Adding coherent and adheres to the PRAM consistency model.
- Memory coherent:
 - for any memory location all processes agree on the same order of all **WRITE** operations to that location.
 - The WRITE operations on different memory location can be in different orders.



Consistency Models – Processor Consistency Model

- How about this sequence?



Consistency Models – Processor Consistency Model

Node N₁

Node N₂

Node N₃

```
{
  ...
  a=d
  c=a+b
  a=4
  ...
}
```

```
{
  ...
  a=10
  ...
  print(a)
}
```

```
{
  ...
  e=foo()
  f=bar()
  ...
}
```

```
...
r1(d)
...
w1(a)
r1(a)
r1(b)
...
w1(c)
w1(a)
...
```

```
...
w2(a)
...
r2(a)
```

Global clock

```
...
w3(e)
...
w3(f)
...
```

Node N₁

Node N₂

Node N₃

Processor Consistency Model

Memory coherences: any memory location all processes agree on the same order of all write operations to that location.

```
...
r1(d)
w3(e)
w3(f)
w2(a)
w1(a)
r1(a)
r1(b)
w1(c)
r1(b)
w1(a)
r2(a)
...
```

```
...
r1(d)
w2(a)
w1(a)
r1(a)
r1(b)
w1(c)
w1(a)
w3(e)
w3(f)
r2(a)
...
```

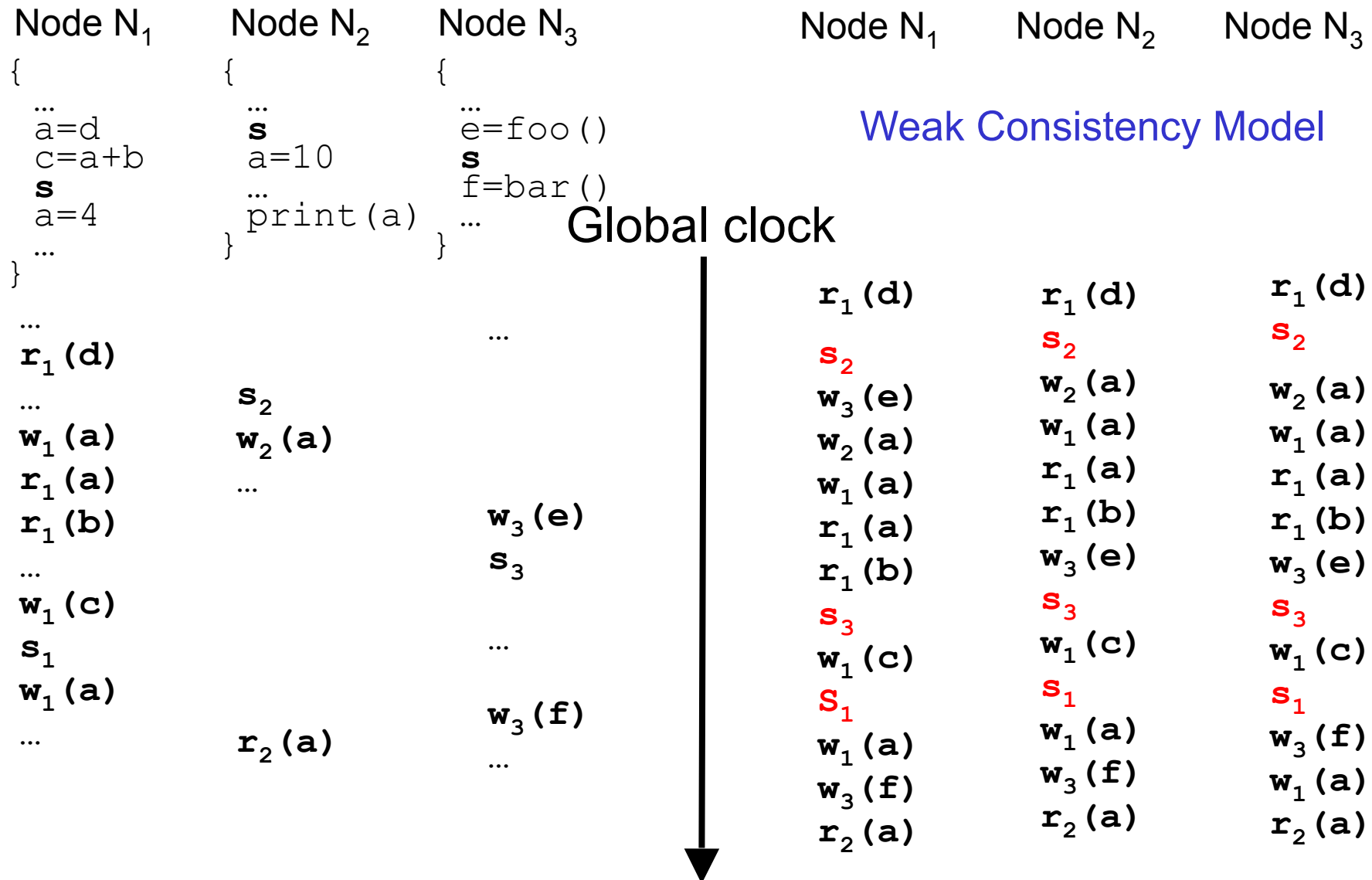
```
...
r1(d)
w2(a)
w1(a)
r1(a)
r1(b)
w1(c)
w1(a)
r2(a)
w3(e)
w3(f)
...
```



Consistency Models – Weak Consistency Model

- Observations by Dubois et al. [1988]:
 - Not necessary to show the change done by every write operation.
 - Isolated access to shared variables are rare.
- Better performance can be achieved if consistency is enforced on a **group** of memory reference operations rather than on **individual** memory reference operations.
- A **synchronization variable** is used to propagate all writes to other machines, and to perform local updates with regard to changes to global data that occurred elsewhere in the distributed system.
- The properties of weak consistency:
 - Accesses to synchronization variables are sequentially consistent.
 - No access to a synchronization variable is allowed to be performed until all previous writes have been completed everywhere. -> To propagate the write before end.
 - No data access (read or write) is allowed to be performed until all previous accesses to synchronization variables have been performed. -> To accept all the updates before start.

Consistency Models – Weak Consistency Model



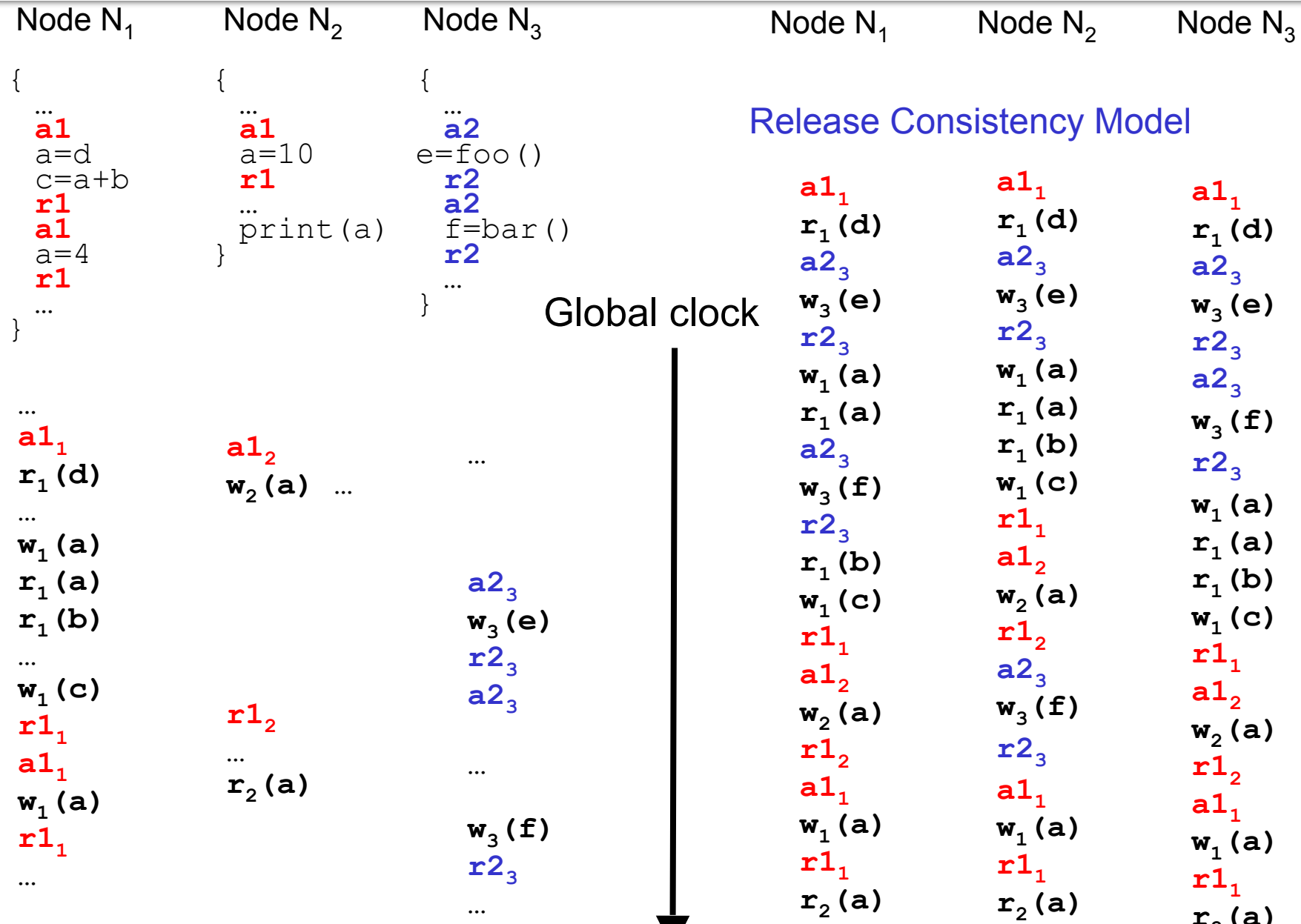
Consistency Models – Release Consistency Model

- Are all the propagations necessary?
 - All changes made to the memory by the process are propagated to other nodes.
 - All changes made to the memory by other processes are propagated from other nodes to the process's node.
- Release consistency mode [Gharachorloo et al. 1990] provides a mechanism to clearly tell the system to decide and perform one of these two operations.
- Two synchronization variables are required:
 - Acquire: a process is about to enter the critical section.
 - Release: a process is about to leave the critical section.
- Programmers are responsible for putting acquire and release at suitable places in their programs.

Consistency Models – Release Consistency Model

- Requirements for release consistency model:
 - All accesses to *acquire* and *release* synchronization variables obey **processor consistency semantics**.
 - All previous *acquires* performed by a process must be completed successfully before the process is allowed to perform a data access operation on the memory.
 - All previous data access operations performed by a process must be completed successfully before a *release* access done by the process is allowed.

Consistency Models – Release Consistency Model



Discussion on Consistency Model

- Which model is most intuitive to you?
- Which model is almost not possible to implement?
- Which model is most intuitive to parallel programming model?
- What are the trade-off for weaker consistent model?

Facebook.com

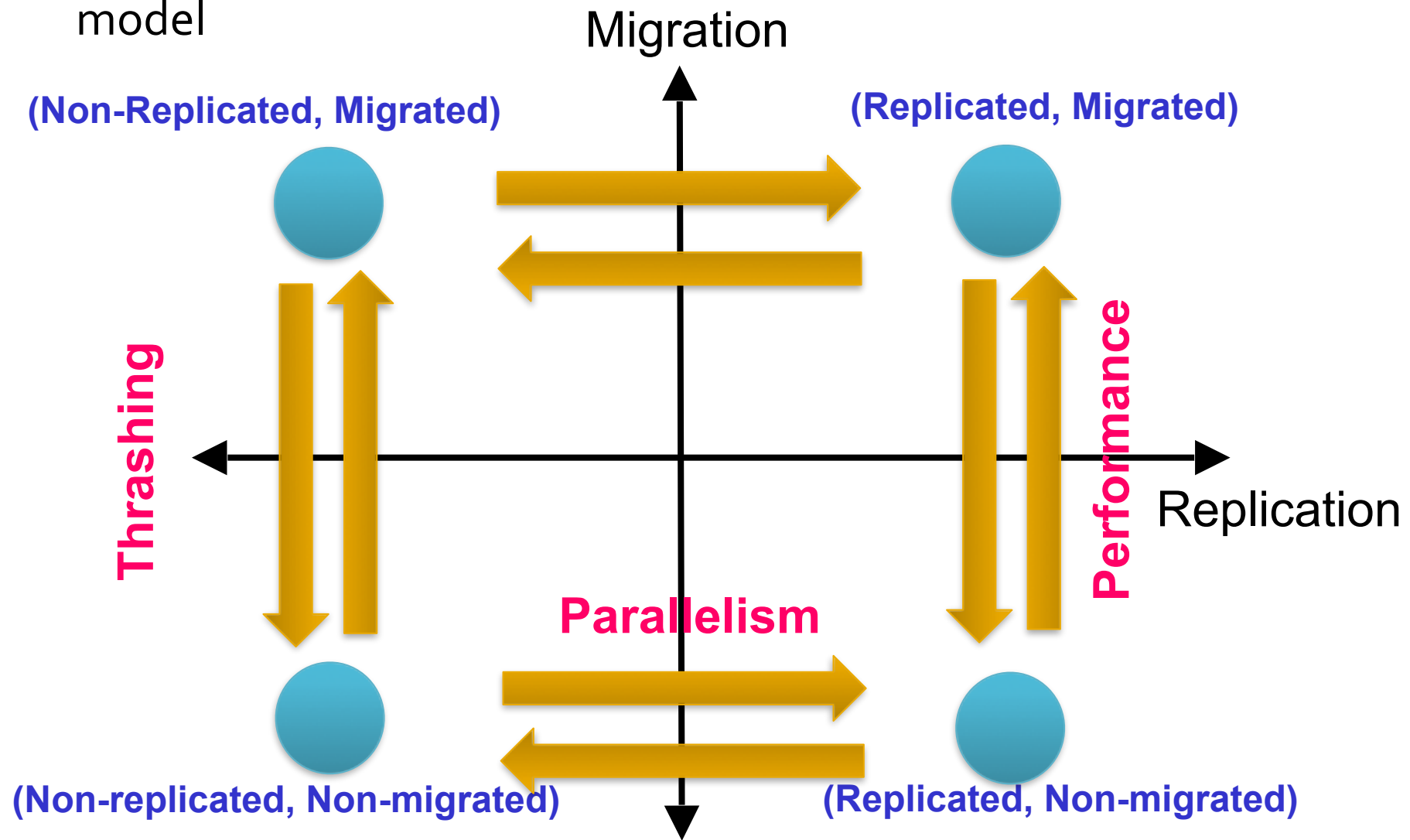
- Suppose facebook.com uses a distributed shared-memory to implement the wall comment/display. Which consistency model should be used so as to minimize the implementation and run-time overhead?

Google Doc

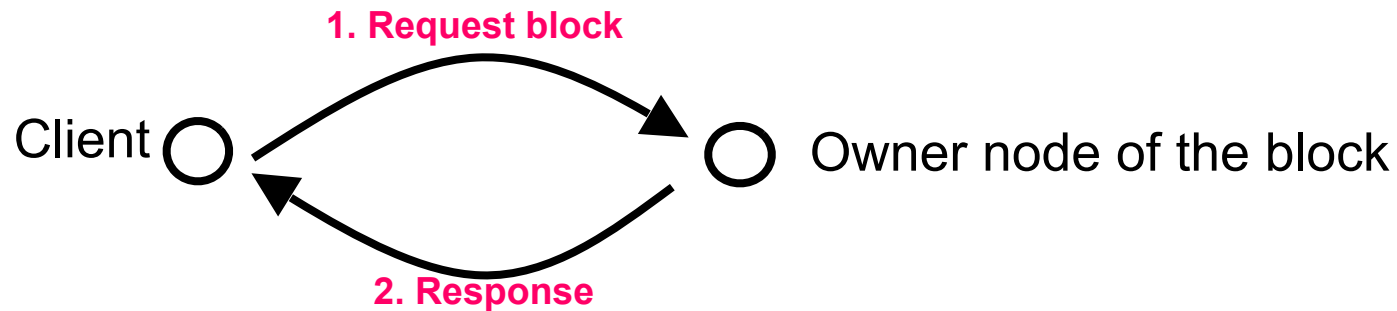
- Suppose that you share your google documents with several groups of friends.
- Which consistency model should be used?
 - One document can be edited by at most one user.
 - One document can be edited by more than one user.
 - 'Save' button is required to store data.
 - No 'Save' button is required to store data.

Implement Sequential Consistency Model

- Not practical to implement strict DSM model.
- Replication and migration strategies for sequential consistency model



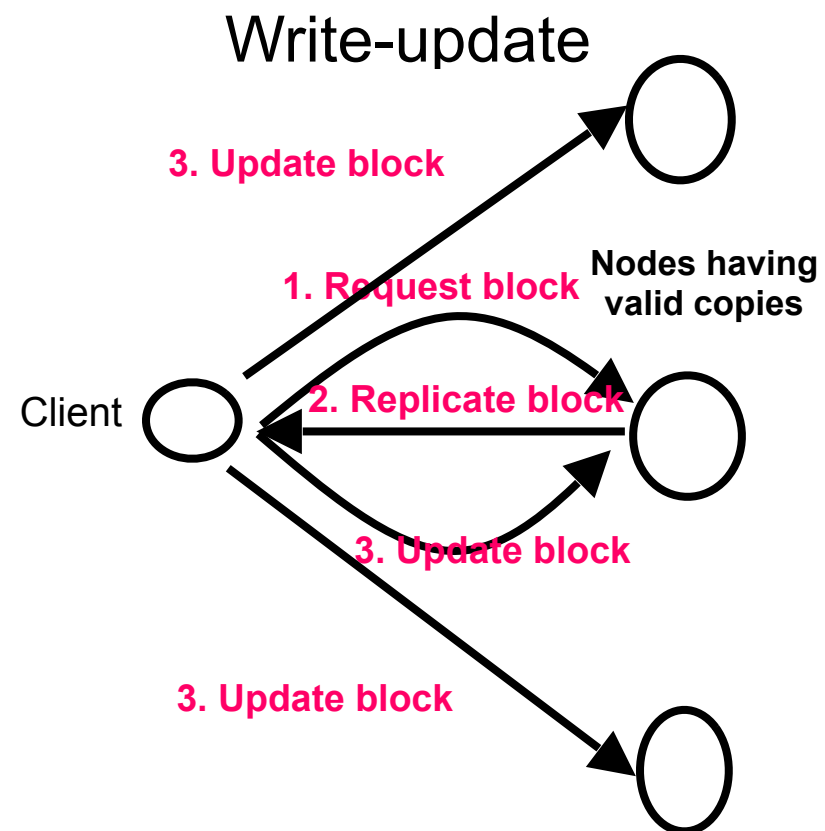
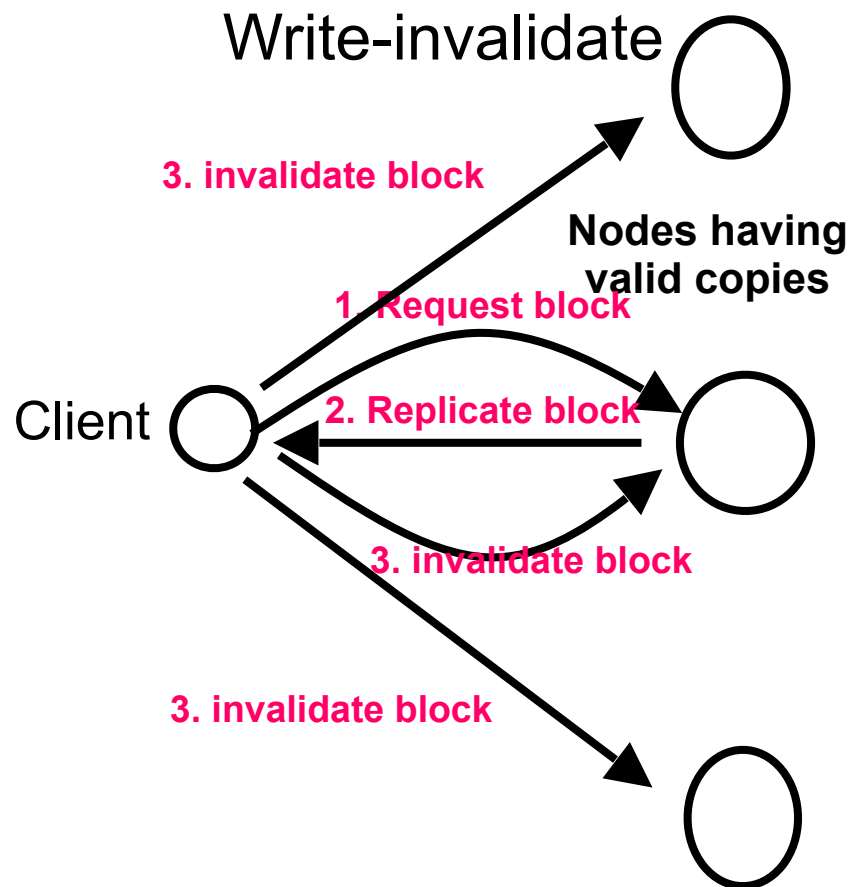
Non-replicated and Non-Migrating Blocks (NRNMB)



- NRNMB strategy:
 - There is a single copy of each block in the entire system.
 - The location of a block never changes.
- NRNMB is easy to implement but has poor performance when the network latency is high.

Replicated and migrated blocks

- Replication complicates the memory coherence protocol.
- Two protocols to ensure sequential consistency.

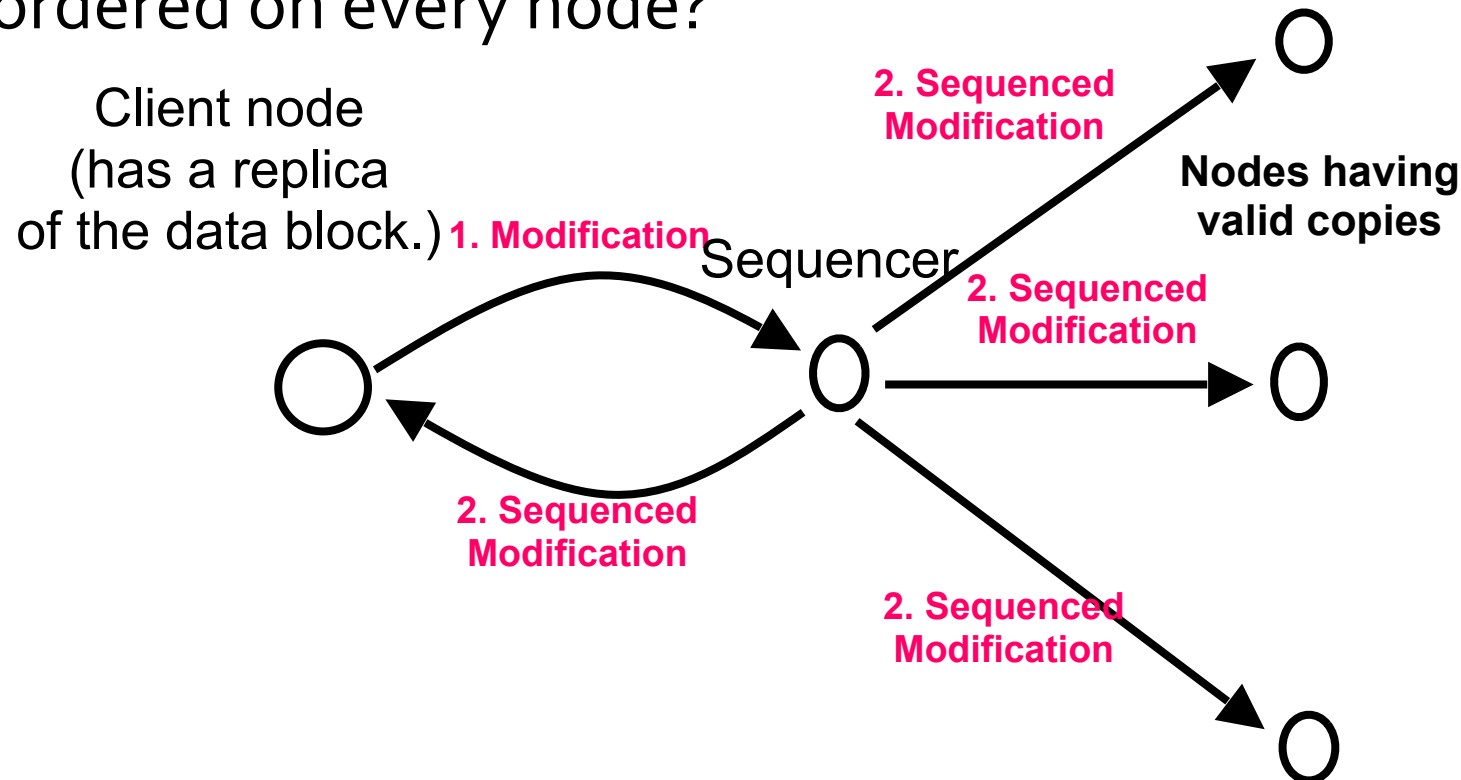


Status Tags for Write-Invalidate Strategy

- The tag indicates
 - whether the block is valid,
 - whether the block is shared, and
 - whether the block is read-only or writeable.
- Read Request
 - If the block is locally available and is valid, the request is satisfied by accessing the local copy.
 - Otherwise, the fault handler generates a read fault and obtains a copy from other nodes.
- Write Request
 - If the block is locally available and is valid and writable, the request is satisfied by accessing the local copy.
 - Otherwise, a fault is generated to obtain a valid copy of the block and changes its status to writable. The fault also invalidates all other copies of the block. Then, the request can be continued.

Global Sequencing Mechanism

- How to assure that the write operations are totally ordered on every node?

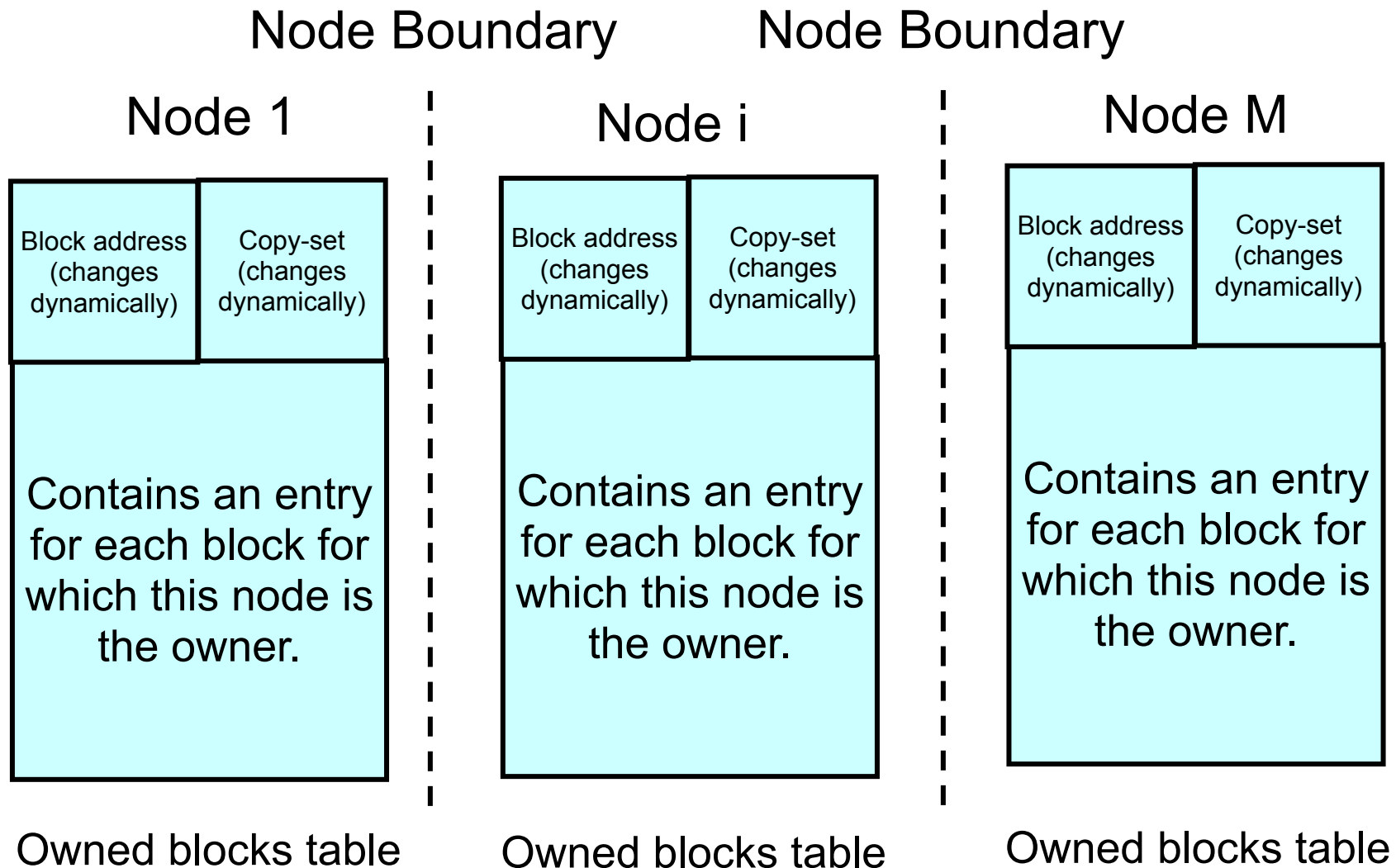


- Virtual clock proposed by Lamport is another approach.
- Write-update is very expensive for use with loosely coupled distributed-memory systems.

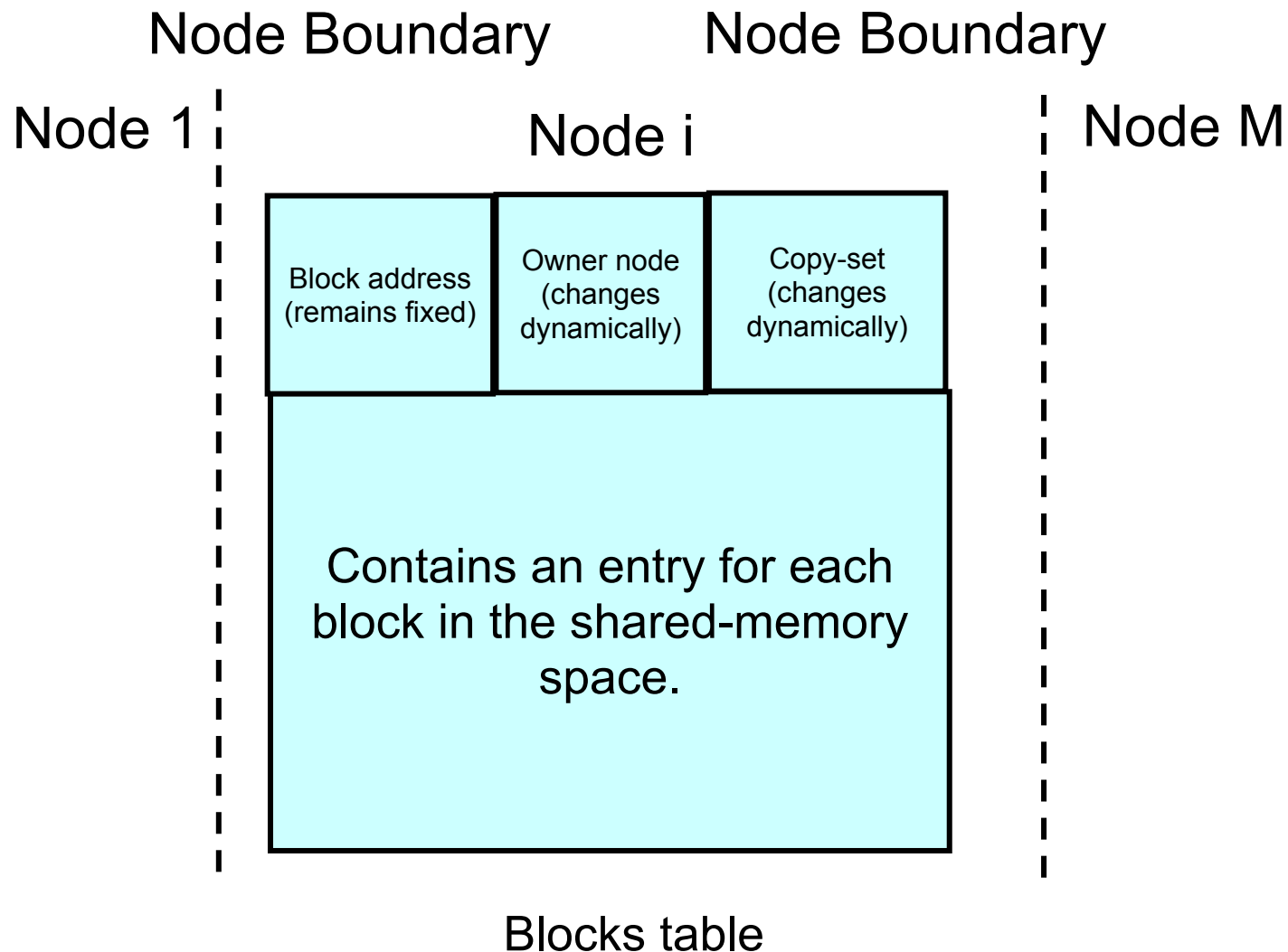
Data Locating in the RMB Strategy

- Data locating issues:
 - Locating the owner of a block.
 - Keeping track of the nodes that currently have a valid copy of the block.
- Possible solutions:
 - Broadcasting
 - Centralized-server algorithm
 - Fixed distributed-server algorithm
 - Dynamic distributed-server algorithm

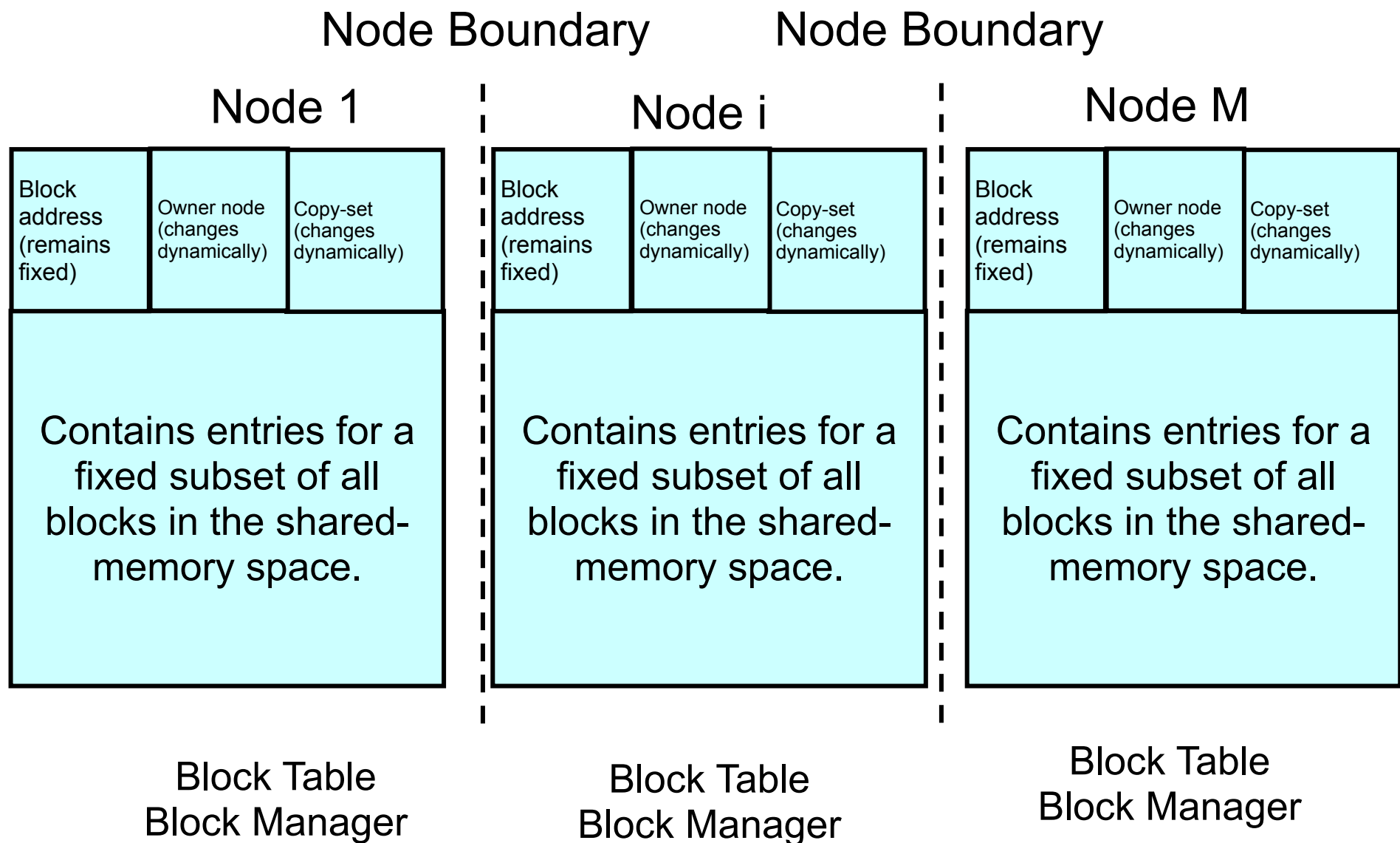
Broadcasting Data Locating Mechanism for RMB Strategy



Centralized-Server Data Locating Mechanism for RMB Strategy



Distributed-Server Data Locating Mechanism for RMB Strategy



Dynamic Distributed-Server Data Locating Mechanism for RMB Strategy

Node Boundary

Node Boundary

Node 1

Node i

Node M

Block address (remains fixed)	Probable Owner node (changes dynamically)	Copy-set (changes dynamically)
Contains entries for a each block in the shared-memory space.		An entry has a value in this filed only if this node is the true owner of the corresponding black.

Block address (remains fixed)	Probable Owner node (changes dynamically)	Copy-set (changes dynamically)
Contains entries for a each block in the shared-memory space.		An entry has a value in this filed only if this node is the true owner of the corresponding black.

Block address (remains fixed)	Probable Owner node (changes dynamically)	Copy-set (changes dynamically)
Contains entries for a each block in the shared-memory space.		An entry has a value in this filed only if this node is the true owner of the corresponding black.

Block Table
Block Manager

Block Table
Block Manager

Block Table
Block Manager

Replacement Strategy

- Challenging Issues for caching shared data:
 - Which block to replace?
 - Where to place a replaced block?
- Replacement Algorithms:
 - Usage-based versus non-usage based: LRU vs. FIFO
 - Fixed space versus variable space
 - Is variable space suitable?

DSM in IVY [Li 1986, 1988]

- Most DSM differentiate the status of data items and use a priority mechanism.
- Each memory block is classified into one of the following five types: unused, nil, read-only, read-owned, and writable.
- Replacement Priority:
 - Both unused and nil have the highest replacement priority. (Note: LRU may leave nil blocks as they are invalidated recently.)
 - Read-only blocks are the next.
 - Read-owned and writable blocks for which replica(s) exist on some other node(s) are the next.
 - Read-owned and writable blocks for which only this node has

Where to place a replaced block

- Two commonly used approaches:
 - Using secondary storage
 - Using the memory space of other nodes.

Thrashing

- Why thrashing?
 - Data blocks are moved back and forth in quick succession.
 - Blocks with read-only permissions are repeatedly invalidated soon after they are replicated.
- Thrashing indicates poor (node) locality in references.
- Avoid thrashing:
 - Providing application-controlled locks
 - Nailing a block to a node for a minimum amount of time
 - Tailoring the coherence algorithms to the shared-data usage patterns