



Sistemas Distribuidos

Módulo 6

**Consistencia, Replicación y Memoria
Compartida Distribuida**



Agenda

1. Replicación
2. Consistencia
 1. Tipos
 2. Modelos
 3. Problemas
3. Replicación y Actualización
 1. Protocolos
4. Memoria Compartida Distribuida



Agenda

1. Replicación
2. Consistencia
 1. Tipos
 2. Modelos
 3. Problemas
3. Replicación y Actualización
 1. Protocolos
4. Memoria Compartida Distribuida



Razones para la Replicación

Hay dos razones principales para la replicación de datos:

- **CONFIABILIDAD**

- Continuidad de trabajo ante caída de la réplica
- Mejor protección contra la corrupción de datos

- **RENDIMIENTO**

- Extenderse en número y en área geográfica (disminuye el tiempo de acceso al dato)
- Consulta simultánea de los mismos datos
- Partición de la red

Precio a pagar por la replicación de datos:
Problemas de Consistencia

Replicación como Técnica de Escalabilidad

Como técnicas para facilitar la expansión se utiliza la replicación y el caching.

Ubicar copias de datos u objetos cercanos a los procesos que los usan mejora el rendimiento por la reducción del tiempo de acceso y resuelve el problema de expansión.

PROBLEMAS:

- La actualización consume más ancho de banda de la red.
- Mantener las copias consistentes resulta un serio problema de escalabilidad y mas en un contexto de consistencia estricta.
- La actualización como una única operación atómica. Se necesitan sincronizar todas las réplicas.

Dilema - Por un lado la replicación tiende a resolver el problema de la expansión (aumenta el rendimiento); por otro mantener consistentes las copias requiere sincronización global.



Agenda

1. Replicación

2. Consistencia

1. Tipos

2. Modelos

3. Problemas

3. Replicación y Actualización

1. Protocolos

4. Memoria Compartida Distribuida

Modelos de Consistencia

Un **MODELO DE CONSISTENCIA** es esencialmente un contrato entre procesos y el almacenamiento de datos.

Es decir: si los procesos acuerdan obedecer ciertas reglas, el almacenamiento promete trabajar correctamente.

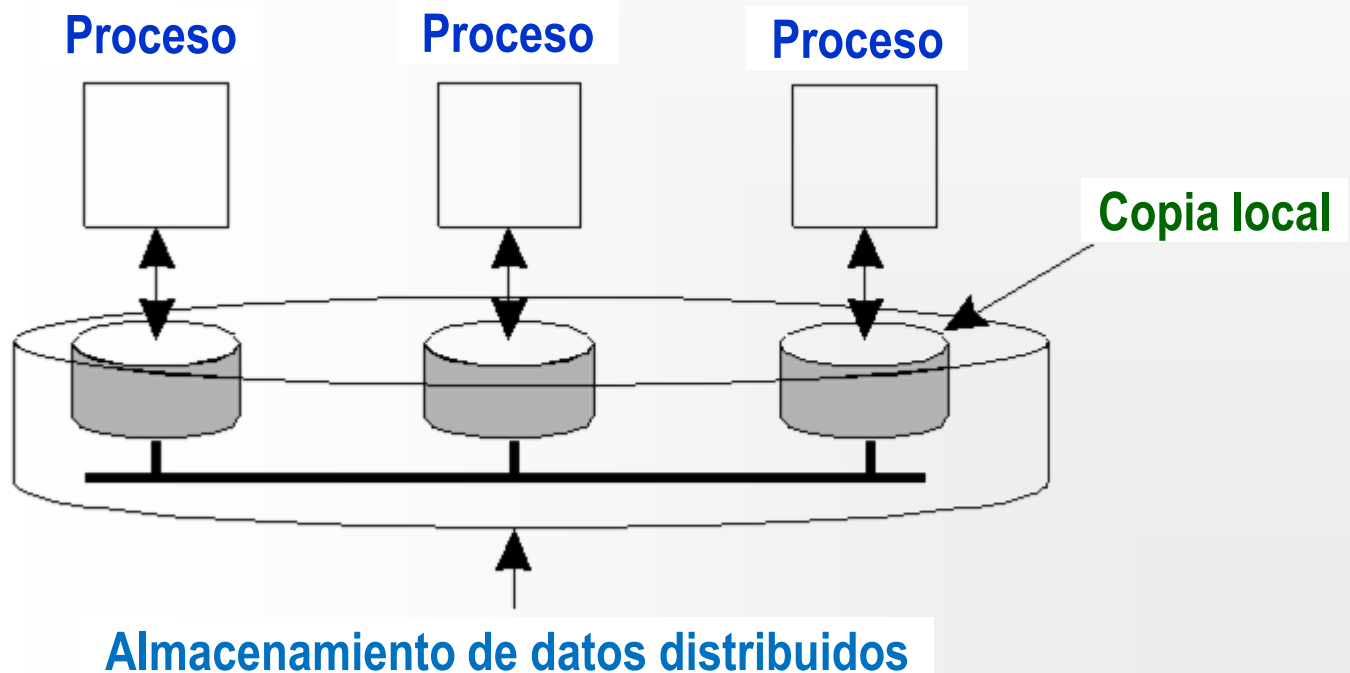
Normalmente un proceso que realiza una operación de lectura espera que esa operación devuelva un valor que refleje el resultado de la última operación de escritura sobre el dato.

Los modelos de consistencia se presentan divididos en dos conjuntos:

- Modelos de consistencia centrados en los datos.
- Modelos de consistencia centrados en el cliente.

Modelos de Consistencia Centrados en Datos

Organización general de un almacenamiento lógico de datos, físicamente distribuidos y replicados a través de múltiples procesos.



Consistencia Estricta

El modelo de consistencia más restrictivo es llamado **consistencia estricta** y es definido por la siguiente condición:

Cualquier lectura sobre un item de dato x retorna un valor correspondiente con la más reciente escritura sobre x

P1: $W(x)a$
P2: $R(x)a$

a) Un almacenamiento estrictamente consistente.

b) Un almacenamiento que no es estrictamente consistente.

P1: $W(x)a$
P2: $R(x)NIL$ $R(x)a$

(b)

La definición supone un ***tiempo global absoluto***

Consistencia - Problema

Solución de Peterson para 2 procesos

Proceso P_i

Datos compartidos

```
int turno;  
boolean flag[2]; inicializado en false
```

repeat

```
flag [i] := true;  
turno := j;  
while (flag [j] and turno = j) do no-op;
```

sección crítica

```
flag [i] := false;
```

sección resto

until *false*;

Consistencia - Problema

// Proc 0

flag[0] = TRUE;

turno = 1;

**while (turno==1 && flag[1]==TRUE)
{};**

// Proc 1

flag[1] = TRUE;

turno = 0;

**while (turno==0 && flag[0]==TRUE)
{};**

Consistencia - Problema

// Proc 0

flag[0] = TRUE;

turno = 1;

while (turno==1 && flag[1]==TRUE)
{

// como flag[1] == FALSE,

// Proc 0 ingresa a la sección crítica

// Proc 1

turno = 0;

flag[1] = TRUE;

while (turno==0 && flag[0]==TRUE)
{

// como turno==1,

// Proc 1 ingresa a la sección crítica

reordenan

Proc 0 observa las dos escrituras desde el Proc 1 fuera de orden, y por esto ingresa a la s.c. antes de observar **flag[i] = true**.

Esto habilita ambos procesos a ingresar a la s.c al mismo tiempo.

Consistencia Secuencial

La **consistencia secuencial** es una forma ligeramente más débil de la consistencia estricta. Satisface la siguiente condición:

El resultado de una ejecución es el mismo si las operaciones (lectura y escritura) de todos los procesos sobre el dato fueron ejecutadas en algún orden secuencial y las operaciones de cada proceso individual aparecen en esta secuencia en el orden especificado por su programa.

P1:	W(x)a		
P2:	W(x)b		
P3:		R(x)b	R(x)a
P4:		R(x)b	R(x)a

P1:	W(x)a		
P2:	W(x)b		
P3:		R(x)b	R(x)a
P4:		R(x)a	R(x)b

- a) Un dato almacenado secuencialmente consistente. (a)
- b) Un dato almacenado que no es secuencialmente consistente. (b)



Consistencia Causal

El modelo de consistencia causal (Hutto and Ahamad, 1990) es un debilitamiento de la consistencia secuencial. Se hace una diferenciación entre eventos que están potencialmente relacionados en forma causal y aquellos que no. Las operaciones que no están causalmente relacionadas se dicen **concurrentes**.

La condición a cumplir para que los datos sean causalmente consistentes es:

Escrituras que están potencialmente relacionadas en forma causal deben ser vistas por todos los procesos en el mismo orden. Escrituras concurrentes pueden ser vistas en un orden diferente sobre diferentes máquinas.

Consistencia Causal

- Esta secuencia es permitida con un almacenamiento causalmente consistente, pero no con un almacenamiento secuencialmente consistente o con un almacenamiento consistente en forma estricta.

P1:	W(x)a		W(x)c	
P2:		R(x)a	W(x)b	
P3:		R(x)a		R(x)c
P4:		R(x)a		R(x)b

P1:	W(x)a			
P2:		R(x)a	W(x)b	
P3:			R(x)b	R(x)a
P4:			R(x)a	R(x)b

(a)

P1:	W(x)a			
P2:			W(x)b	
P3:			R(x)b	R(x)a
P4:			R(x)a	R(x)b

(b)

(a) Una violación de una consistencia causal.

(b) Una correcta secuencia de eventos en una consistencia causal.

Consistencia First Input First Output

- Condición Necesaria

Escrituras realizadas por un proceso único son vistas por los otros procesos en el orden en que son hechas, pero escrituras desde diferentes procesos pueden ser vistas en diferente orden por diferentes procesos.

P1:	W(x)a			
P2:	R(x)a	W(x)b	W(x)c	
P3:			R(x)b	R(x)a
P4:			R(x)a	R(x)b

- Una secuencia válida de eventos de una consistencia First Input First Output.

Consistencias – Agrupación de Operaciones

Las consistencias son modeladas a través de mecanismos de sincronización para exclusión mutua y transacciones.

- Entrar_SC() – asegura que todos los datos están actualizados en su almacenamiento local para realizar las operaciones necesarias.
- Salir_SC()



Sección Crítica

Idea Básica

No importa que las lecturas y escrituras de una serie de operaciones sean inmediatamente conocido por otros procesos. Solo se requiere que el efecto de la serie sea conocido.



Consistencia y Coherencia

- Un **modelo de consistencia** describe lo que se puede esperar con respecto a ese **conjunto** cuando múltiples procesos operan simultáneamente sobre **esos datos**.
- Un **modelo de coherencia** describe lo que se puede esperar para **un solo elemento de dato**. Se asume que cada elemento de dato está replicado.



Modelos de Consistencia Centrados en el Cliente

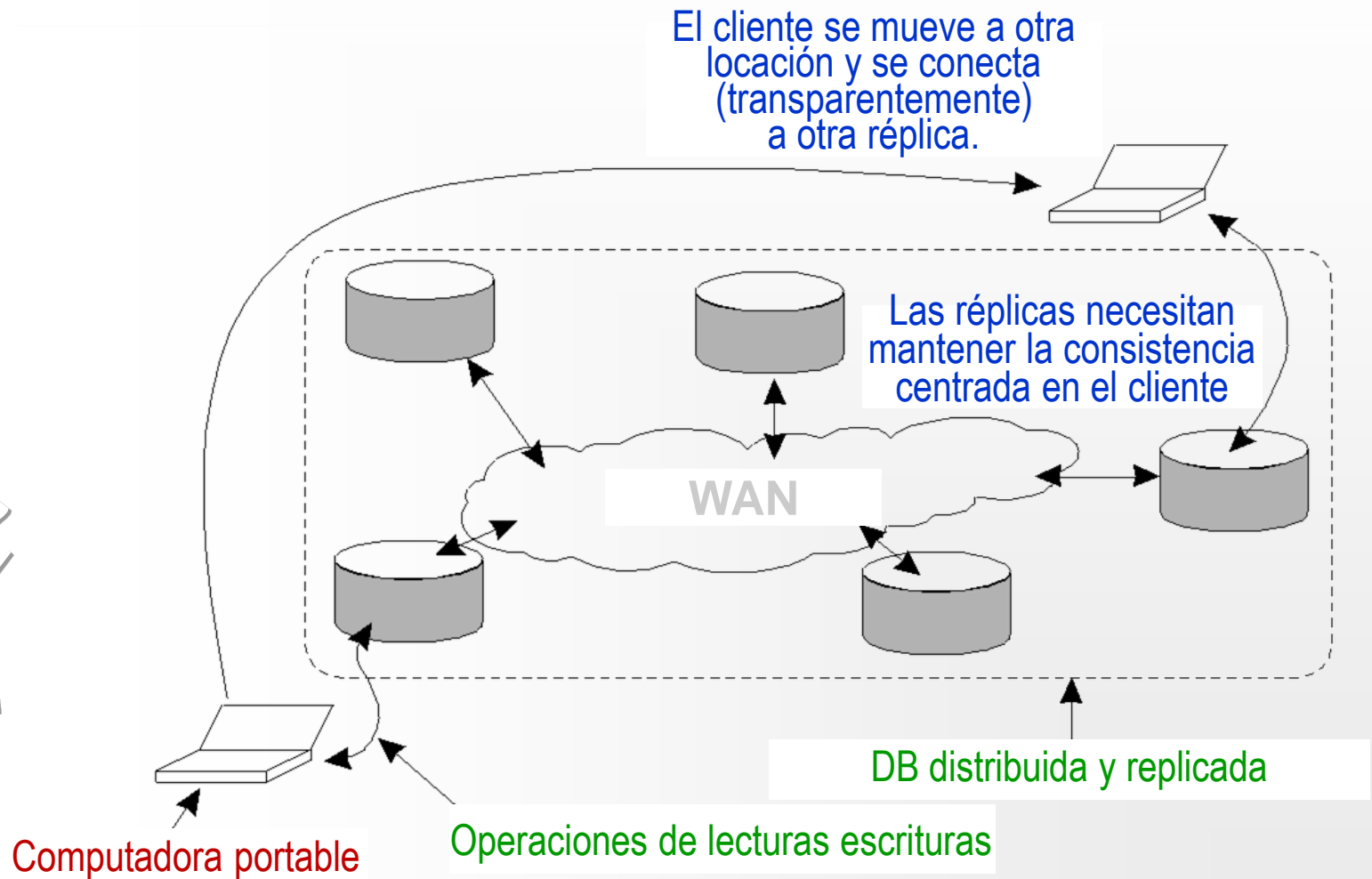
Clase especial de almacenamiento de datos distribuidos.

- Están caracterizados por la falta de actualizaciones simultáneas, o cuando dichas actualizaciones ocurren, pueden ser fácilmente resueltas.
- La mayoría de las operaciones son de lectura.

La introducción de **modelos de consistencia centrados en el cliente** permiten esconder muchas inconsistencias de manera relativamente fácil.

- En esencia, provee garantías **para un único cliente** concerniente a la consistencia de accesos a los datos de ese cliente.
- No se dan garantías para accesos concurrentes por diferentes clientes.

Consistencia Centrada en el Cliente





Consistencias Centradas en el cliente

- Lecturas Monotónicas
- Escrituras Monotónicas
- Lea sus Escrituras
- Escrituras seguidas de Lecturas

Referencias:

- [1] Distributed Systems: Principles and Paradigms, 1ra ed, Tanenbaum et ál, items 6.3.2, 6.3.3, 6.3.4 y 6.3.5
- [2] Distributed Systems, 3ra ed, Van Maarten y Tanenbaum, items 7.3



Agenda

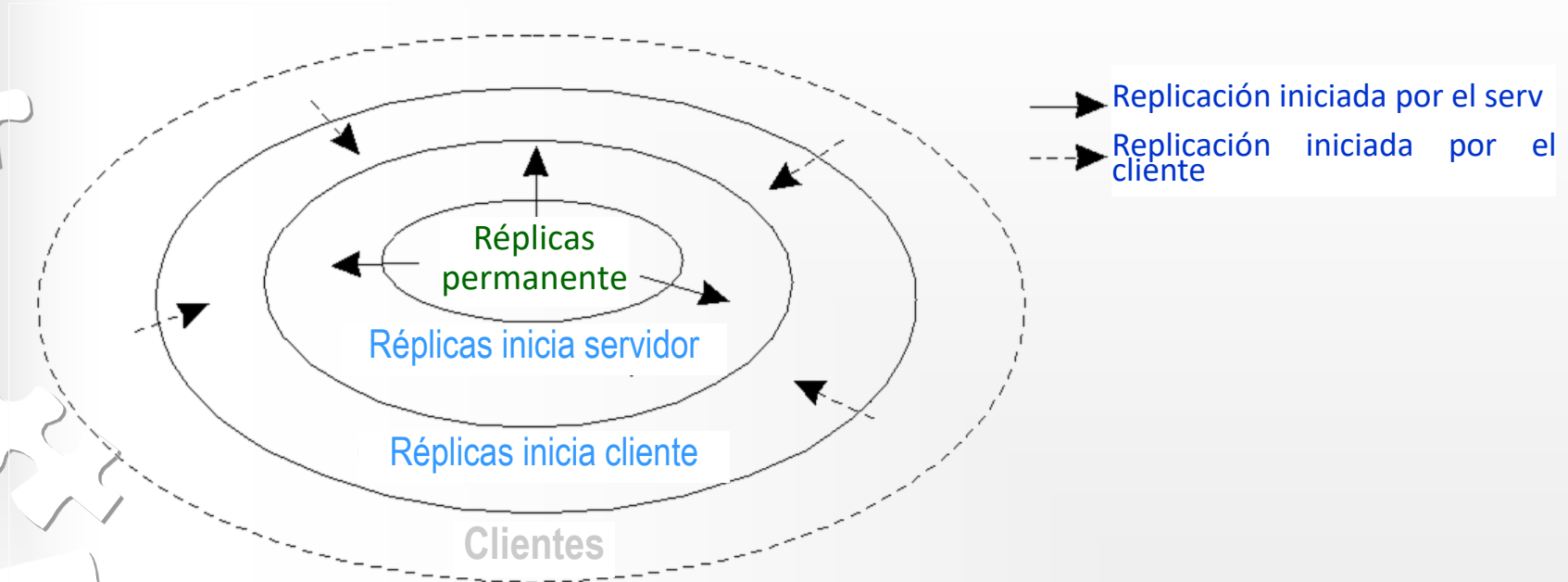
1. Replicación
2. Consistencia
 1. Tipos
 2. Modelos
 3. Problemas
3. Replicación y Actualización
 1. Protocolos
4. Memoria Compartida Distribuida



Replicación: Ubicación del Contenido

- **RÉPLICAS PERMANENTES:** proceso / máquina siempre teniendo una réplica
- **RÉPLICAS INICIADAS POR EL SERVIDOR:** proceso que puede alojar dinámicamente una réplica en solicitud de otro servidor en el almacén de datos
- **RÉPLICAS INICIADAS POR EL CLIENTE:** proceso que puede alojar dinámicamente una réplica en solicitud de un cliente (caché del cliente)

Ubicación de Réplicas



La organización lógica de diferentes clases de copias de datos almacenados en tres anillos concéntricos

Actualización de Réplicas

ESTADO VERSUS OPERACIONES – QUÉ SE PROPAGA

1. Propagar solamente la notificación de una actualización.



PROTOCOLOS DE INVALIDACIÓN

2. Transferir datos de una copia a otra.

3. Propagar la operación de actualización hacia otras copias.



REPLICACIÓN ACTIVA



Actualización de Réplicas

ACTUALIZACIONES SE INSERTAN O SE EXTRAEN?

- Método basado en push – protocolos basados en servidor
- Método basado en pull – protocolos basados en cliente



Protocolos de Consistencia

Un *protocolo de consistencia* describe una implementación de un modelo específico de consistencia.

Los modelos de consistencia en los cuales las operaciones están globalmente seriadas son los modelos más importantes y los más ampliamente aplicados.

Estos modelos incluyen consistencia secuencial, consistencia débil con variables de sincronización y transacciones atómicas.

Los protocolos pueden ser:

- Protocolos basados en el primario
- Protocolos de escritura replicada



Protocolos Basados en el Primario

Cada item de dato x tiene en el almacenamiento de datos un primario asociado, el cual es responsable de coordinar las operaciones de escritura sobre x .

Se debe hacer una distinción si el primario está fijo en un servidor remoto o si las operaciones de escritura pueden ser llevados localmente luego de mover el primario a donde reside el proceso que inició la operación de escritura.



Protocolos Basados en el Primario

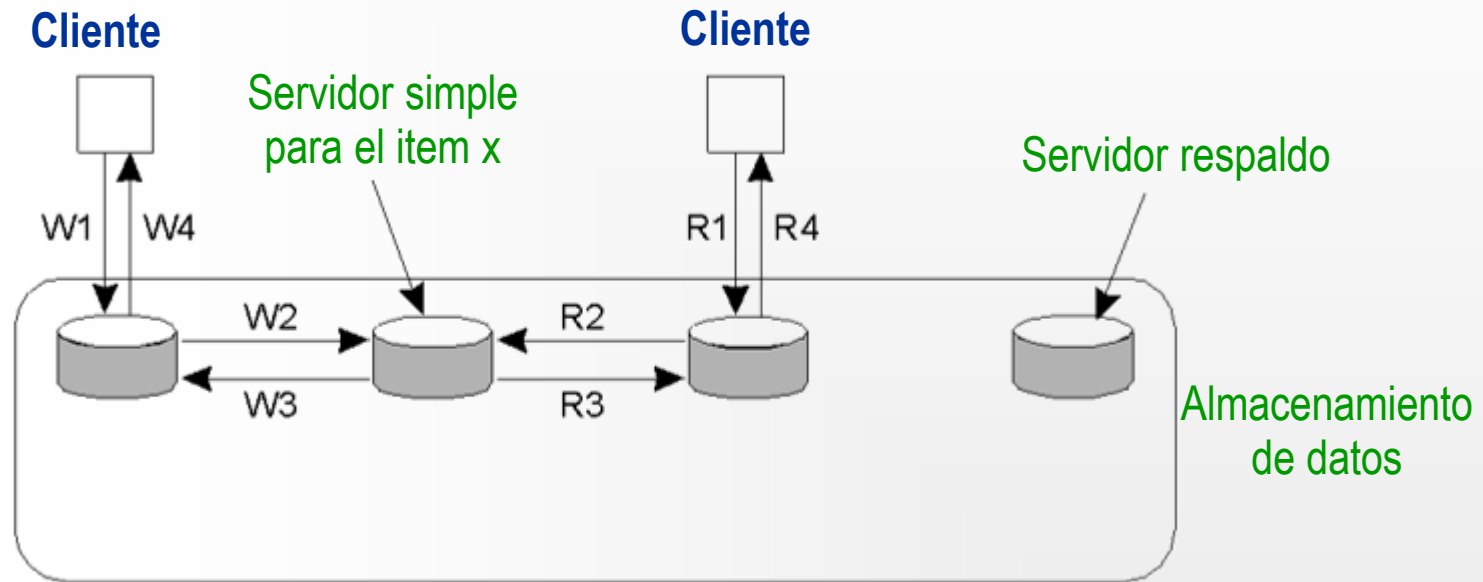
Cada **item de dato x** tiene en el almacenamiento de datos un **primario asociado**, el cual es responsable de coordinar las operaciones de escritura sobre x .

Formas de implementación

- El primario está fijo en un servidor remoto
- El primario se mueve. Las operaciones de escritura pueden ser realizadas localmente luego de mover el primario a donde reside el proceso que inició la operación de escritura.

Protocolos Escritura Remota

Protocolo de escritura remota basado en primario con **servidor fijo** al cual las lecturas y escrituras son encaminadas.



W1. Requerimiento Escritura

W2. Req. encaminado al serv. por x

W3. ACK escritura completa

W4. ACK escritura completa

R1. Requerimiento lectura

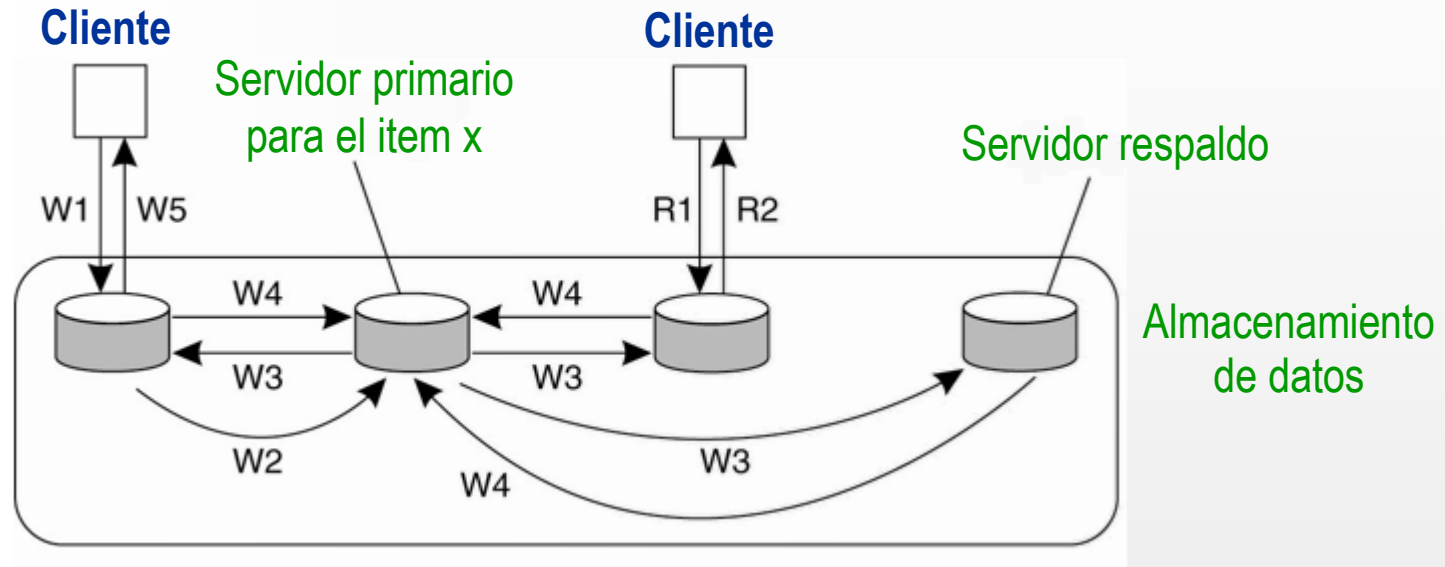
R2. Req. encaminado al serv. por x

R3. Retorno respuesta

R4. Retorno respuesta

Protocolos Escritura Remota

El principio de protocolo de respaldo primario.



W1. Solicitud escritura

W2. Mueve el requerimiento al primario

W3. Actualiza los respaldos

W4. ACK actualización

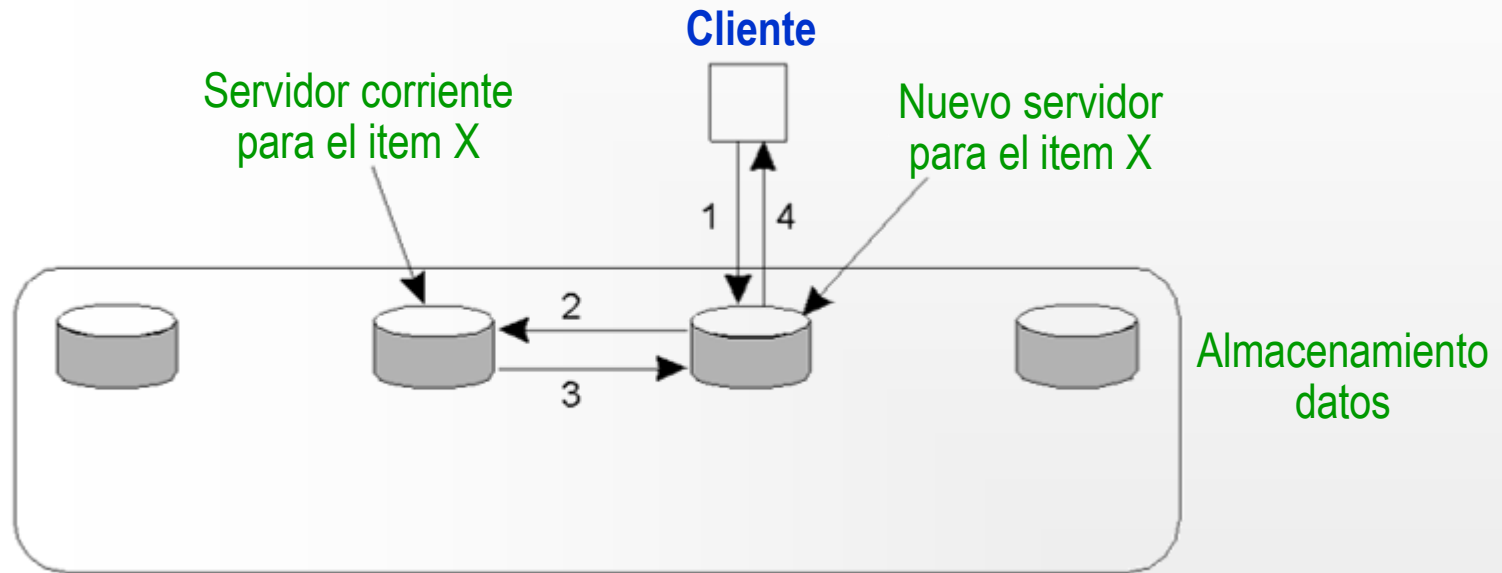
W5. ACK escritura completada

R1. Solicitud de lectura

R2. Respuesta a la lectura

Protocolos de Escrituras Locales

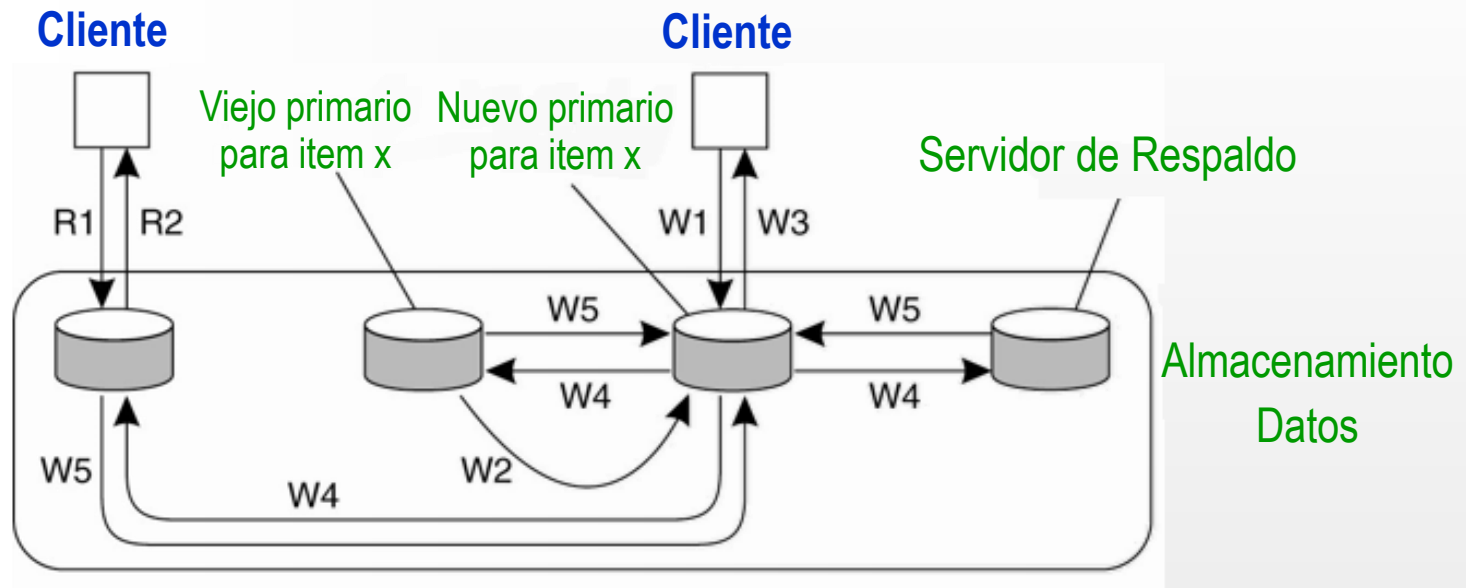
Protocolo escritura local basado en primario en el cual una única copia es migrada entre procesos.



1. Solicitud de lectura o escritura.
2. Continua solicitud al servidor corriente de x.
3. Se mueve el ítem x al servidor del cliente.
4. Retorna el resultado de la operación sobre el servidor del cliente.

Protocolos de Escritura Local

Protocolo de respaldo primario en el cual el primario migra hacia el proceso que espera realizar una actualización.



W1. Solicitud escritura

W2. Mueve ítem x a un nuevo primario

W3. ACK escritura completada

W4. Actualiza los respaldos

W5. ACK actualización

R1. Solicitud de lectura

R2. Respuesta a la lectura

Protocolos de Escritura Replicada

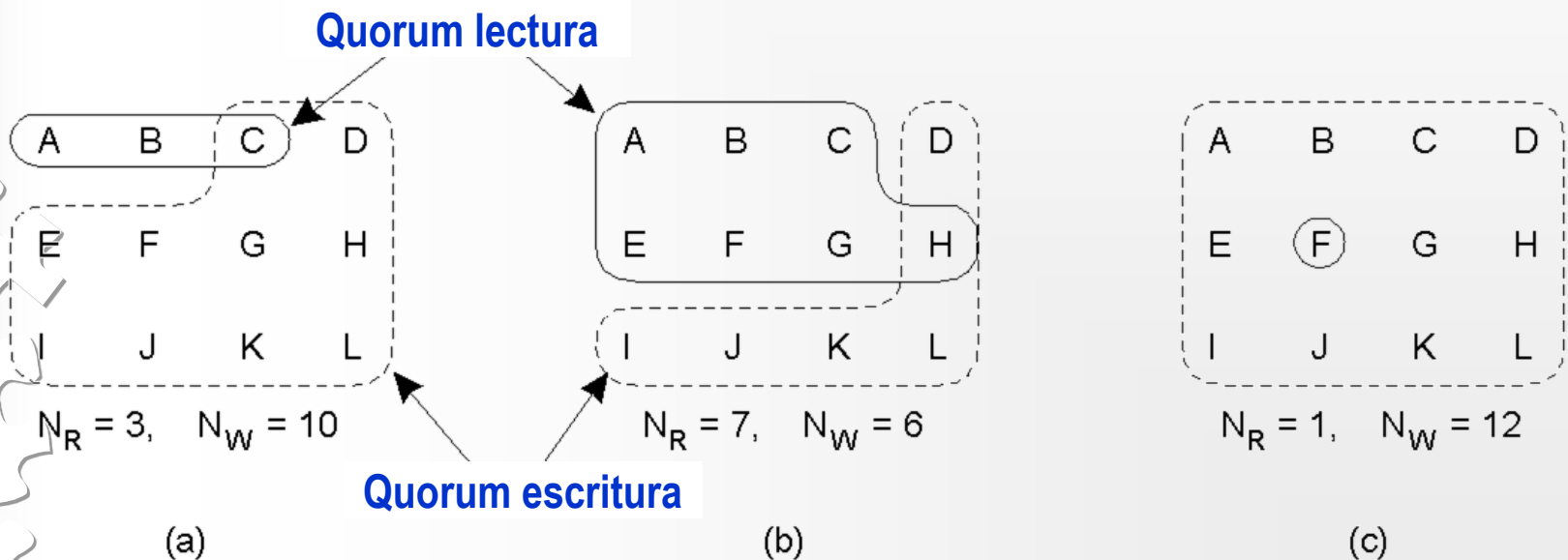
En los protocolos de escrituras replicadas, las operaciones pueden ser llevadas a cabo sobre múltiples réplicas en vez de una sola como en el caso de las réplicas basadas en el primario.

- Replicación Activa
- Protocolos basados en quorum
 - Restricciones de los quorums (lectura y escritura)
 - $N_R + N_W > N$
 - $N_W > N/2$

Protocolos Basados en Quorum

Tres ejemplos del algoritmo de votación:

- a) Una correcta elección de conjuntos de lectura y escritura.
- b) Una elección que puede llevar a conflictos escritura-escritura.
- c) Una elección correcta, conocida como ROWA (read one, write all)





Agenda

1. Replicación
2. Consistencia
 1. Tipos
 2. Modelos
 3. Problemas
3. Replicación y Actualización
 1. Protocolos
4. Memoria Compartida Distribuida

Memoria Compartida Distribuida

Comunicación entre procesos

- **Pasaje de Mensajes**
- **Memoria Compartida**

Memoria Compartida Distribuida (MCD)

La capa abstracta se implementa en:

Kernel

Rutinas en tiempo de corrida



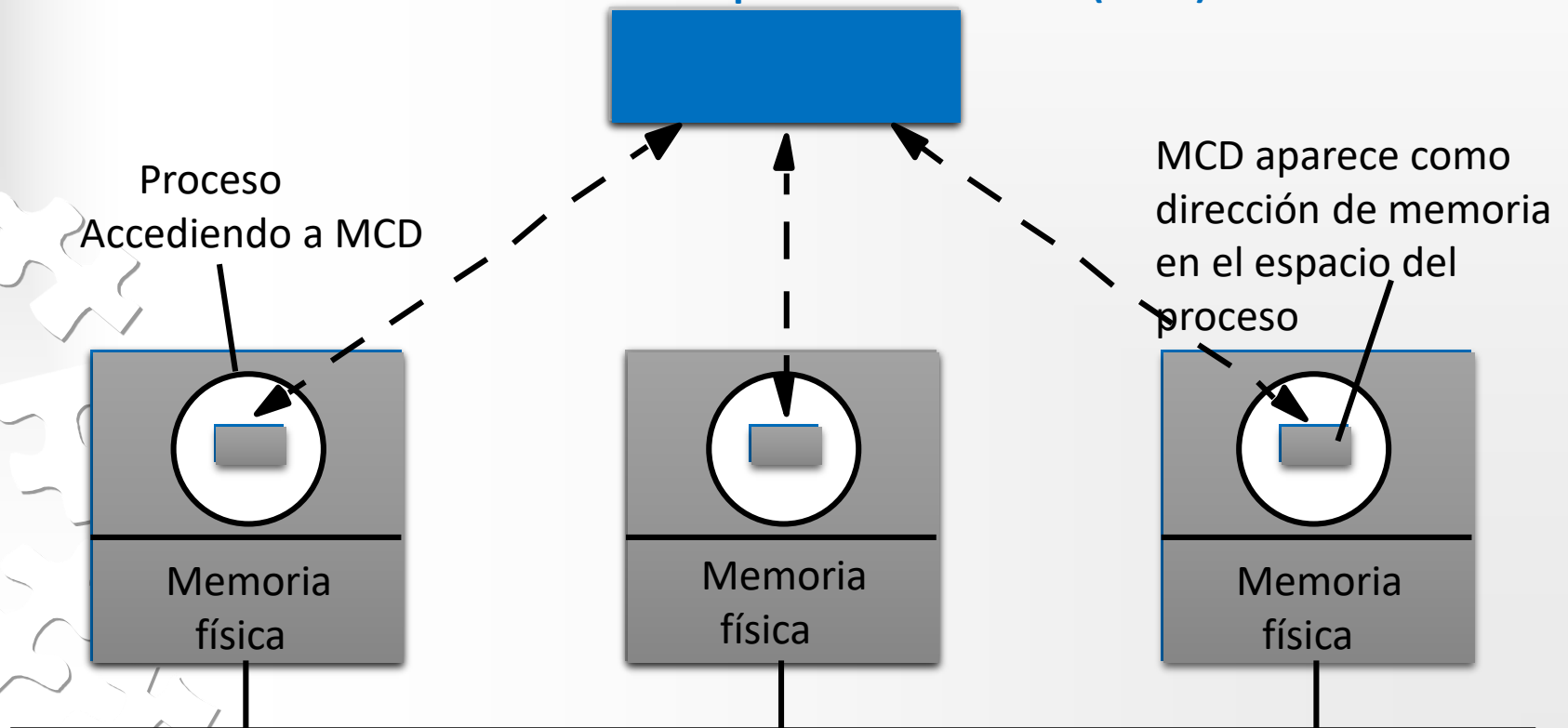
**Dónde es
apropiado?**

Memoria Compartida Distribuida

La memoria se parte en “bloques” fijos. “*Caching*” evita la *latencia de acceso*.

Puede soportar: **migración y/o replicación**

Memoria Compartida Distribuida (MCD)





Memoria Compartida Distribuida

DISEÑO E IMPLEMENTACIÓN DE MCD - ASPECTOS

- Granularidad
- Estructura del espacio de memoria compartida
- Coherencia de memoria y sincronización de acceso
- Localización de datos y accesos
- Estrategias de reemplazo
- Thrashing
- Heterogeneidad



Memoria Compartida Distribuida - GRANULARIDAD

Problemas que se presentan con el tamaño del bloque

- Sobrecarga de paginado
- Tamaño de directorio
- Thrashing
- Falso compartir



Memoria Compartida Distribuida

Ventajas de usar un tamaño de bloque igual a la página de memoria local:

- Permite el uso de sistemas existentes de falta de páginas
- Permite el control de derechos de acceso
- Si las páginas pueden colocarse en un *paquete* no impone sobrecarga en la red
- El tamaño es adecuado con respecto a la contención de memoria



Memoria Compartida Distribuida - ESTRUCTURA DEL ESPACIO

La estructura define la abstracción de la vista a los programadores de aplicaciones en un sistema MCD.

Puede ser vista para unos como palabra y para otros como objetos. La estructura y la granularidad está relacionadas.

No estructurada: es un arreglo de palabras.

Por tipo de datos: colección de objetos (Clouds y Orca) o colección variables en lenguaje fuente (Munin y Midway). La granularidad es la variable o el objeto.

Como base de datos: implica una memoria asociativa.



Memoria Compartida Distribuida – MODELO DE CONSISTENCIA: SECUENCIAL

Ejemplo de Implementación

- Es realizable
- Soporta semántica más intuitiva para la coherencia de memoria.
- No impone nada extra al programador.

Desventaja: baja concurrencia

La consistencia débil y de liberación con sus variables implícitas proveen mejor concurrencia y soportan una semántica más intuitiva.



Memoria Compartida Distribuida - MODELO DE CONSISTENCIA: SECUENCIAL - IMPLEMENTACIÓN

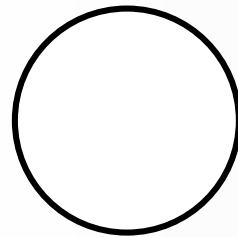
Protocolos: dependen de la estrategia elegida

Estrategias:

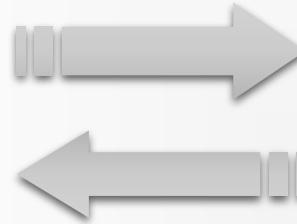
- **NO RÉPLICA, NO MIGRATORIA (NRNMB)**
- **NO RÉPLICA, MIGRATORIA (NRMB)**
- **RÉPLICA, MIGRATORIA (RMB)**
- **RÉPLICA, NO MIGRATORIA (RNMB)**

MCD - MODELO DE CONSISTENCIA: SECUENCIAL – IMPLEMENTACIÓN NRNMB

Nodo Cliente
Requerimiento-Respuesta

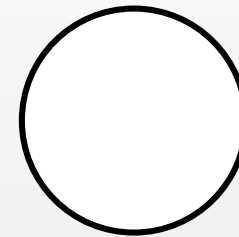


requerimiento



respuesta

Nodo dueño del bloque
Recibe el requerimiento,
hace el acceso, envía
respuesta



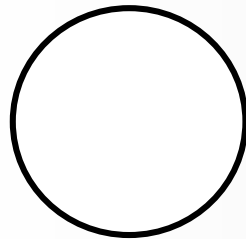
Desventajas:

- La serialización de los requerimientos crea un *cuello de botella*.
- No es posible paralelismo

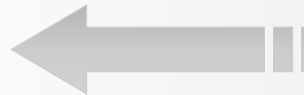
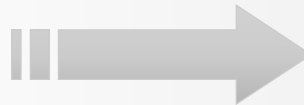
Memoria Compartida Distribuida - NRMB

NO RÉPLICA, MIGRATORIA (NRMB)

Nodo Cliente
(Se convierte en el
nuevo dueño del bloque
luego de su migración)

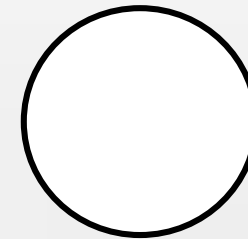


**requerimiento
bloque**



migra bloque

Nodo dueño del bloque
(dueño del bloque antes
de la migración)





Memoria Compartida Distribuida - NRMB

Ventajas:

- No se incurre en costos de comunicación cuando un proceso accede a un dato local
- Permite tomar ventajas de la localidad de accesos a los datos

Desventajas:

- Problemas de *thrashing*
- No hay paralelismo



Memoria Compartida Distribuida - NRMB

Localización de los bloques

1. Broadcasting

- No escala bien. El sistema de comunicaciones es *cuello de botella*.
- La latencia de la red es grande



Memoria Compartida Distribuida - NRMB

2. Algoritmo servidor centralizado

Todos los nodos conocen su dirección.

- Reduce paralelismo (serializa los requerimientos).
- Una falla hace caer el sistema.



Memoria Compartida Distribuida - NRMB

3. Algoritmo servidor distribuido fijo

- Distribuye el rol del servidor.
- Manejan bloques en varios nodos.
- Cada uno maneja un subconjunto determinado de bloques de datos.



Memoria Compartida Distribuida - NRMB

4. Algoritmo servidor distribuido dinámico

- Cada nodo tiene su propia tabla de información de los bloques de la MCD.
- No siempre la información de esta tabla es correcta.



Memoria Compartida Distribuida - RMB

RÉPLICA, MIGRATORIO (RMB)

Ataca la desventaja principal de las estrategias no replicantes: la reducción del paralelismo.

Se incrementa el costo de las operaciones de escritura.

Se debe mantener la **consistencia**

La replicación complica el protocolo de coherencia de memoria.



Memoria Compartida Distribuida - RMB

Protocolos:

- Escritura invalidante
- Escritura actualizada

Este esquema requiere un ordenamiento total de las operaciones de escritura para lograr la consistencia secuencial.

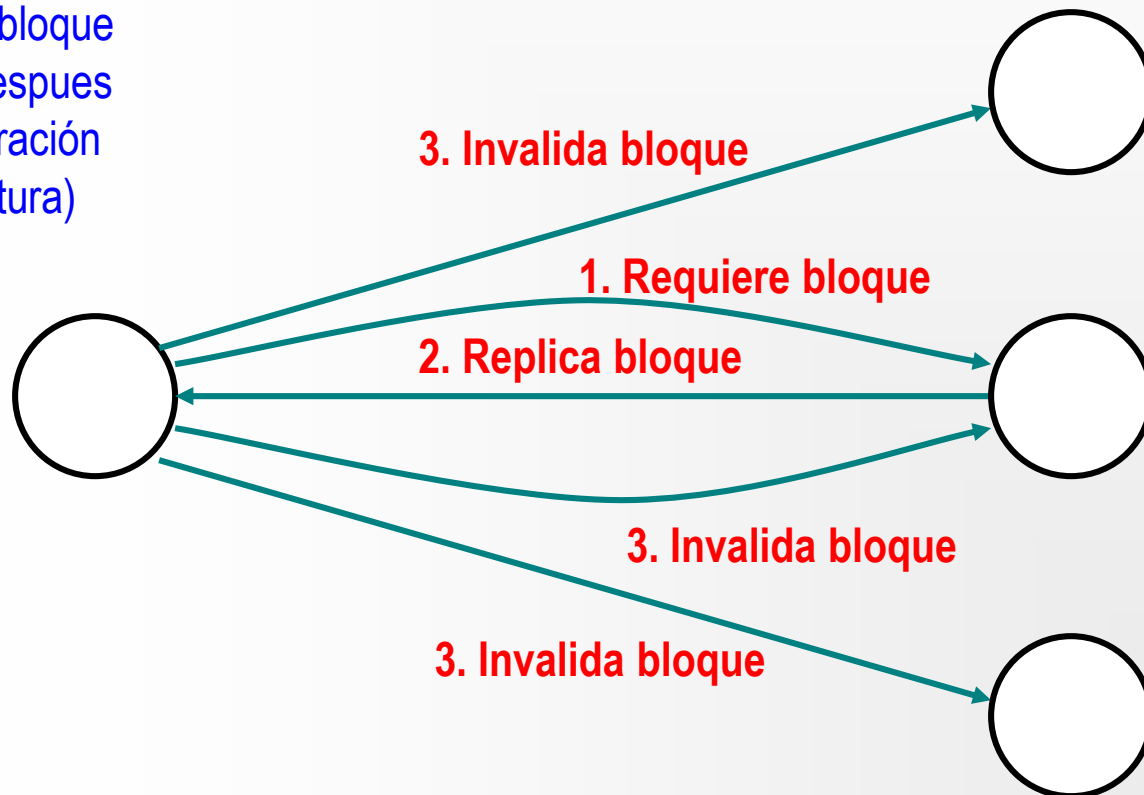
Se utiliza un secuenciador global.

Memoria Compartida Distribuida - RMB

Escritura invalidante

Nodo Cliente
(tiene la copia
válida del bloque
de dato despues
de la operación
de escritura)

Nodos teniendo copias
válidas de bloques de
datos antes de escritura

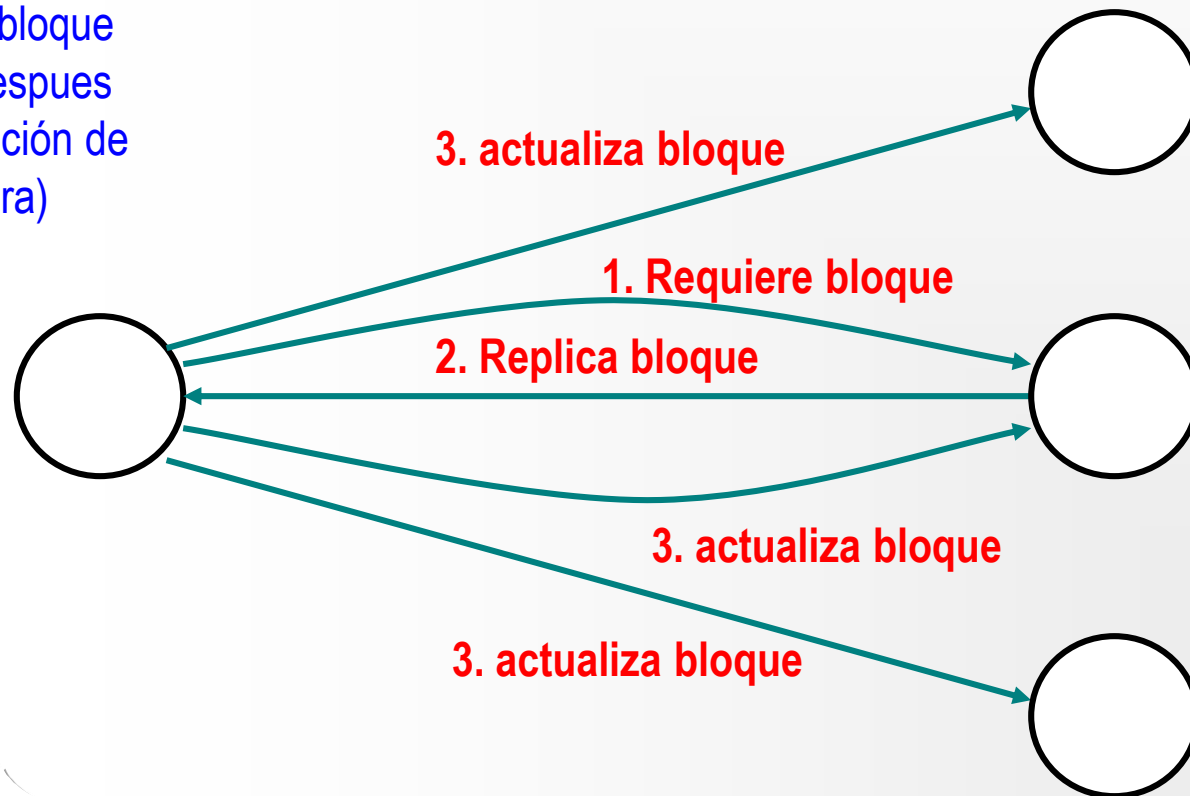


Memoria Compartida Distribuida - RMB

Escritura actualizada

Nodo Cliente
(tiene la copia
válida del bloque
de dato despues
de la operación de
escritura)

Nodos teniendo copias
válidas de bloques de
datos antes de escritura

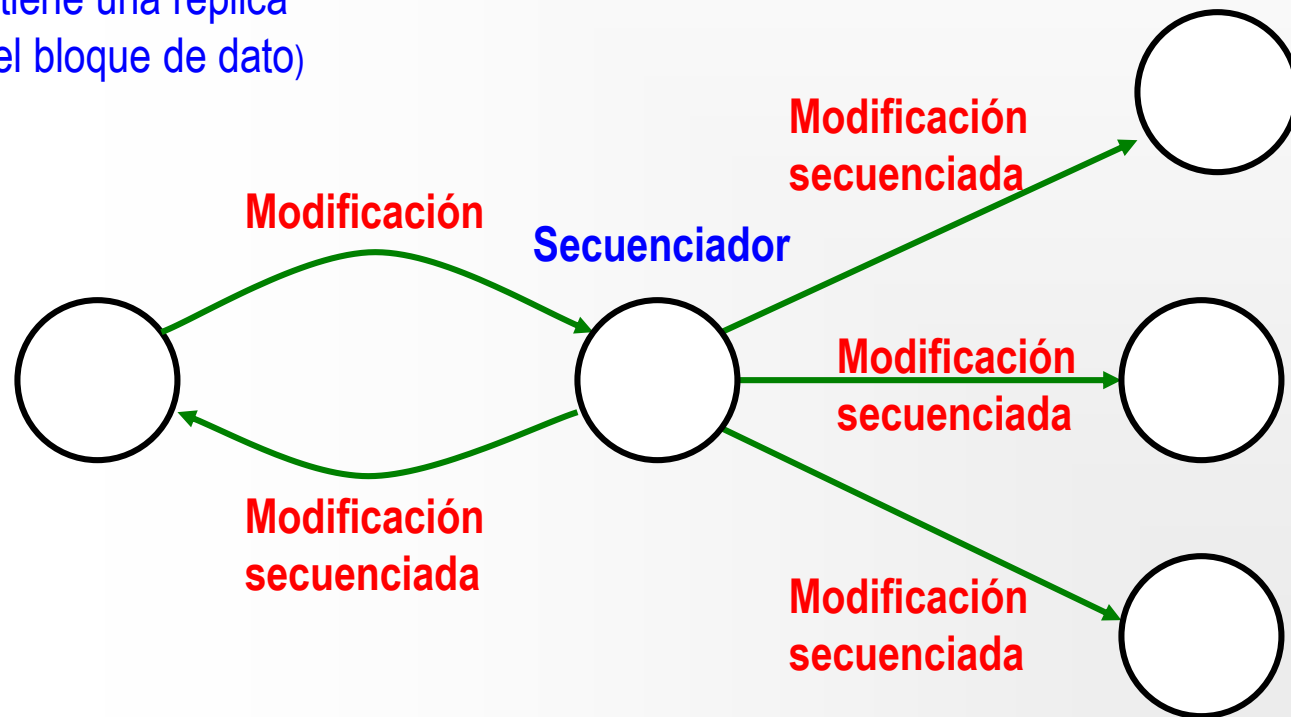


Memoria Compartida Distribuida - RMB

Mecanismo secuenciador

Nodo Cliente
(tiene una réplica
del bloque de dato)

Otros nodos teniendo una
réplica del bloque de datos





Memoria Compartida Distribuida - RMB

Localización de datos

Necesita:

1. Localizar al dueño de un bloque.
2. Tener la pista de los nodos que corrientemente tienen una copia válida del bloque.

Algoritmos usados:

- Broadcasting.
- Algoritmo servidor centralizado.
- Algoritmo servidor distribuido fijo.
- Algoritmo servidor distribuido dinámico.



Memoria Compartida Distribuida - RMB

Para reducir la cadena de nodos a atravesar para alcanzar el verdadero dueño del bloque, la tabla de bloques del nodo es adaptada de la siguiente forma:

- a) Siempre que el nodo recibe un requerimiento de invalidación.
- b) Siempre que el nodo pasa la propiedad.
- c) Siempre que el nodo pasa un requerimiento que falla.



Memoria Compartida Distribuida - RNMB

Réplica, no migratoria (RNMB)

La ubicación de cada réplica es fija.

Siempre que se escribe se actualiza.

Localización de datos

- a) Las localizaciones de las réplicas nunca cambian.
- b) Todas las réplicas se mantienen consistentes.
- c) Sólo un requerimiento de lectura puede ser directamente enviado a uno de los nodos teniendo una réplica del bloque y todos los requerimientos de escritura deben primero ser enviados al secuenciador.



Memoria Compartida Distribuida - RNMB

Estructura:

- ✓ Una tabla de bloques por cada nodo.
- ✓ Una tabla de secuencias en el secuenciador.

La tabla del secuenciador tiene los siguientes campos:

- Dirección del bloque.
- Conjunto de réplicas.
- Número de secuencia



Memoria Compartida Distribuida

Estrategias de reemplazo

- a) Qué bloque tiene que ser reemplazado.
- b) Dónde debe ubicarse el bloque reemplazado.

¿Qué bloque tiene que ser reemplazado?

- a) Usado vs. no usado (LRU).
- b) Espacio fijo vs. espacio variable.

Ubicación el bloque reemplazado

- Usar almacenamiento secundario.
- Usar memoria de otros nodos.



Memoria Compartida Distribuida

En algunos sistemas cada bloque es clasificado:

- No usado
- Nil (invalidado)
- Read-only
- Read-owned: RO y el nodo es dueño
- Escribible

De acuerdo a ésto, la prioridad de reemplazo sería:

1. No usado y Nil
2. Read-only
3. Read-owned
4. Escribible



Memoria Compartida Distribuida

Thrashing

Cuando hay migración y se producen las siguientes situaciones:

- Datos entrelazados entre distintos nodos.
- Bloques con permiso RO son repetidamente invalidados inmediatamente que son replicados.

Implica poca **localidad**



Memoria Compartida Distribuida

Soluciones:

1. Proveer locks de aplicación controlada: fija el bloque por un tiempo.
2. No permitir que quiten, por un tiempo t , un bloque del nodo (t está dado por estadística o por accesos).
3. Ajustar el algoritmo de coherencia para usar modelos de datos compartidos.



Memoria Compartida Distribuida

Ventajas

- Abstracción mas simple.
- Mejor portabilidad de programas de aplicación distribuidos.
- Mejor desempeño de algunas aplicaciones. Lo hace posible la *localidad de los datos*, el *movimiento de datos por demanda* y *espacio de direcciones mas grande*.
- Ambiente de comunicación mas flexible.
- Facilidad para la migración de procesos.
- Comunicación indirecta, asincrónica y persistente.



Bibliografía:

- Sinha, P. K.; “Distributed Operating Systems: Concepts and Design”, IEEE Press, 1997.
- Tanenbaum, A.S.; van Steen, Maarten; “Distributed Systems: Principles and Paradigms”. 2nd Edition, Prentice Hall, 2007 and 1st Edition 2002.
- van Steen, Maarten; Tanenbaum, A.S.; “Distributed Systems”. 3rd Edition, Prentice Hall, 2017.
- Coulouris, G.F.; Dollimore, J. y T. Kindberg; “Distributed Systems: Concepts and Design”. 5th Edition Addison Wesley, 2011.