

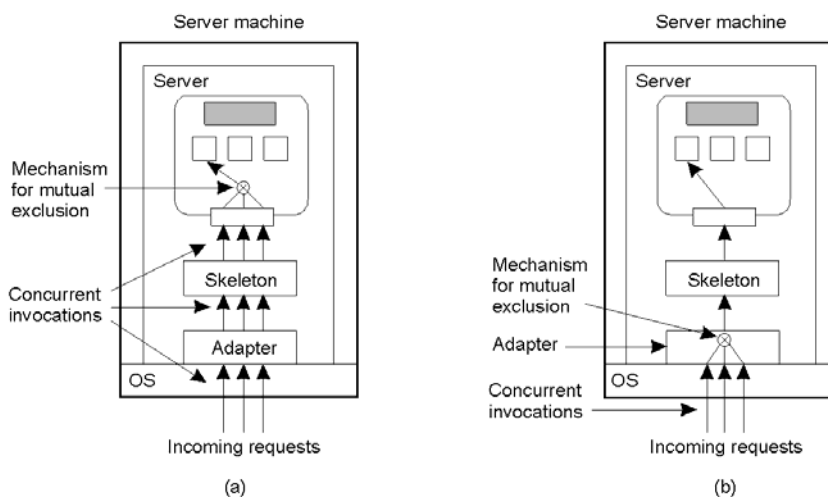
Consistența și replicarea datelor în Internet

Bazat pe "Distributed Systems" de AS Tanenbaum

Motive pentru replicare

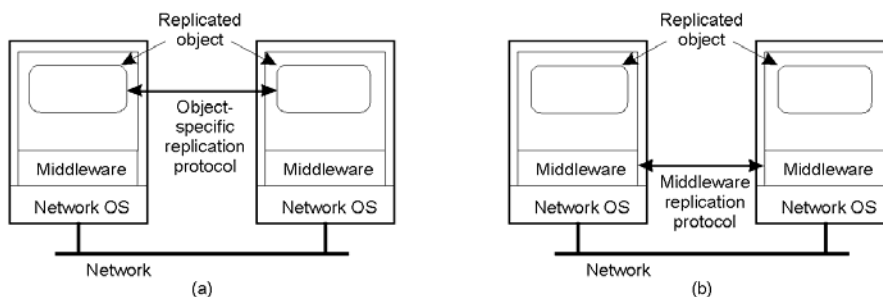
- Siguranta
 - Sistemul continua sa lucreze dupa caderea unei replici
 - Protectie mai buna impotriva datelor corupte
- Performanta
 - Scalare in numar resurse
 - Scalare in arie geografica
- Replicare => Probleme de consistenta

Controlul accesului concurrent la un obiect



- a) Manipularea accesului concurrent de catre obiect
 b) Folosirea unui adaptor pentru manevrarea accesului concurrent

Invocari concurente ale obiectelor replicate



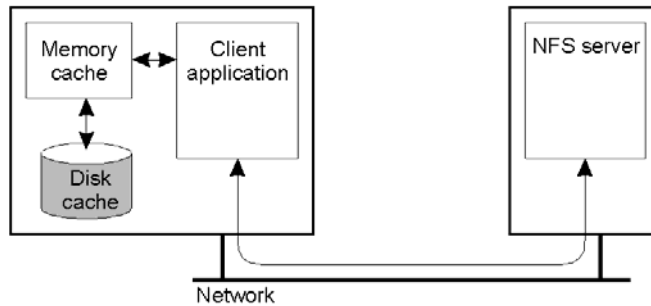
Problema: aceleasi modificari de stare sa apara la obiectele replicate

Solutia: invocarile sa se faca in aceeasi ordine la toate replicile

Doa **implementari** posibile:

- Obiectele sunt "constiente" de replicare (Globe)
- Sistemul distribuit (middleware) face gestiunea replicilor (Piranha)
 - de ex. ordonarea totala a invocarilor la toate replicile

Replicare in sistemele de fisiere



Problema

- eficientizarea accesului partajat la fisiere

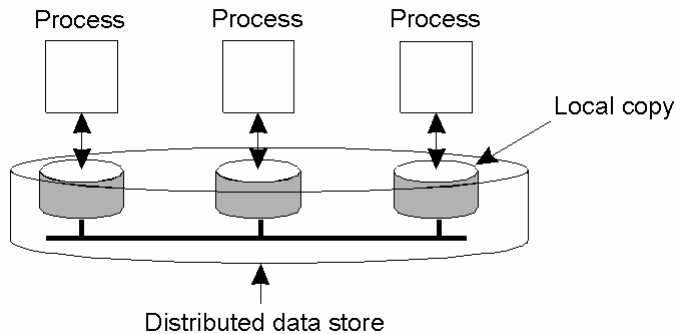
Solutii

- replicare la client - caching pentru date, atribute, handler, directoare
- replicarea fisierelor la server sau in sisteme p2p

Modele de consistență

- Consistenta stricta - greu de pastrat
- Solutie = constrangeri mai slabe de consistenta
- Depind de
 - tiparele de acces si actualizare
 - aplicatii
- Modele de consistenta
 - centrate pe date
 - centrate pe client

Modele de consistenta centrate pe date



Organizare generala depozit de date distribuit si replicat.
Operatii: **read** (copie locala), **write** (propagata la celelalte copii)

Model de consistenta = contract intre proces si data store

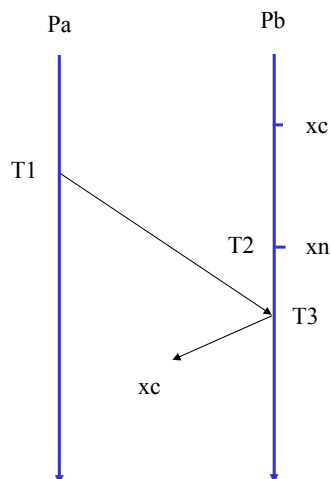
- **read** ar trebui sa intoarca rezultatul ultimei operatii

Consistența Strictă

Orice citire a unei date x intoarce valoarea rezultata din cea mai recenta scriere a lui x .

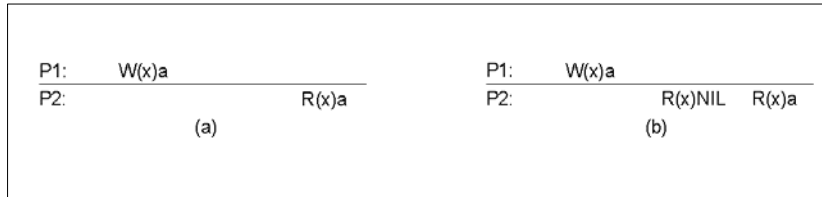
Ex:

- Doua masini A si B
- Doua procese: Pa pe A si Pb pe B
- x memorat pe B (doar); valoarea curenta este x_c
- La T_1 , Pa citește x (trimite un mesaj lui B sa citeasca x)
- La $T_2 > T_1$, Pb scrie x ; noua valoare este x_n
- La $T_3 > T_2$, mesajul de citire de la Pa soseste la B
- Consistenta stricta \Rightarrow B ar trebui sa intoarca x_c



Consistența Strictă (2)

- Se pastreaza ordinea absoluta globala in timp
- Toate scrierile sunt instantaneu vizibile tuturor proceselor



Comportarea a doua procese operand pe acelasi item.

- a) Memorie strict consistenta.
- b) Memorie care nu este strict consistenta.

Greu de implementat

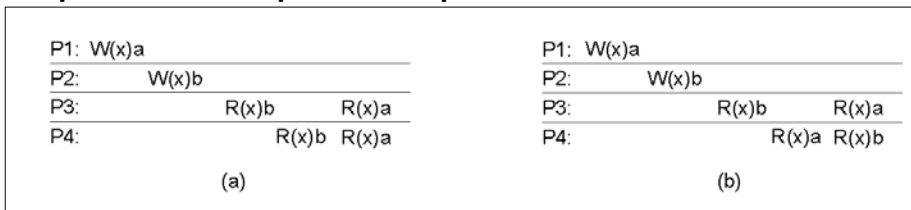
Programele concurente nu sunt bazate pe timpul global sau pe viteza proceselor (ex. Producer / Consumer)

Sunt necesare modele mai slabe

Consistența Secvențială (1)

Consistența secvențială: Rezultatul oricarei execuții este același cu cel obținut dacă operațiile (read, write) tuturor proceselor asupra depozitului de date sunt executate într-o secvență oarecare și operațiile fiecărui proces individual apar în această secvență în ordinea specificată de programul său.

- Orice întretesere validă a operațiilor read și write este acceptabilă
- Toate procesele văd aceeași întretesere de operații
- Operațiile nu au amprente de timp



Un depozit de date **(a)** secvențial consistent și **(b)** care nu este secvențial consistent.

Consistența Secvențială (2)

x, y, si z sunt initializate la 0

Proces P1	Proces P2	Proces P3
x = 1; print (y, z);	y = 1; print (x, z);	z = 1; print (x, y);

Trei procese concurente.

assignment = write
print = doua read simultane

- Sunt posibile 90 ordonari valide diferite ale instructiunilor
90 = 720 (=6!) / 8
- mai putin de 64 tipare de **signatura**
 - signatura = concatenarea iesirilor lui P1, P2 si P3 in aceasta ordine
- procesele trebuie sa accepte toate rezultatele valide

Consistența Secvențială (3)

x = 1; print (y, z); y = 1; print (x, z); z = 1; print (x, y);	x = 1; y = 1; print (x,z); print(y, z); z = 1; print (x, y);	y = 1; z = 1; print (x, y); print (x, z); x = 1; print (y, z);	y = 1; x = 1; z = 1; print (x, z); print (y, z); print (x, y);
Prints: 001011	Prints: 101011	Prints: 010111	Prints: 111111
Signatura: 001011	Signatura: 101011	Signatura: 110101	Signatura: 111111
(a)	(b)	(c)	(d)

Patru secvente de executie valide pentru procesele anterioare.

Axa verticala este timpul.

Nu toate semnaturile sunt permise. Ex. 000000 si 001001

Consistența Cauzală (1)

Operatii:

- Legate Cauzal
 - ex:
 - procesul p executa (1) write x;
 - ulterior, procesul q executa (2) read x; (3) write y
 - (1) si (3) sunt legate causal
- Concurente

Conditie necesara:

Operatiile "write" care sunt potential legate causal trebuie sa fie vazute de toate procesele in aceeași ordine. Scrierile concurente pot fi vazute in diferite ordini pe masini diferite.

Consistența Cauzală (2)

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

$W_1(x)a$ si $W_2(x)b$ sunt dependente causal;

$W_2(x)b$ si $W_1(x)c$ sunt concurente.

Aceasta secventa este permisa cu o memorie causal-consistenta, dar nu cu una sequential sau strict consistenta.

Consistentă Cauzală (3)

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) Violare a consistenței cauzale. (Cele două scrieri sunt legate cauzal.)
- b) O secvență corectă de evenimente într-o memorie cauzal-consistentă. (Cele două scrieri nu mai sunt legate cauzal.)

Consistentă FIFO (1)

Condiție Necesară:

Scriverile făcute de un singur proces sunt văzute de toate celelalte procese în ordinea în care au fost executate, dar scrierile din procese diferite pot fi văzute în ordini diferite de procese diferite.

Usor de implementat

- etichetând fiecare operație `write` cu perechea (proces, număr secvență)

Consistenta FIFO (2)

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

O secventa valida de evenimente pentru consistenta FIFO

Consistenta FIFO (3)

Process P1	Process P2
x = 1;	y = 1;
if (y == 0) kill (P2);	if (x == 0) kill (P1);

Doua procese concurente.

- Consistenta FIFO: ambele procese pot fi omorate
 - scrierile sunt vazute in ordini diferite
- Consistenta Secventiala: sase posibile intreteseri, nici una cu ambele procese omorate

Gruparea operatiilor

Datele partajate pot fi asociate cu o **variabila de sincronizare**.

Un proces foloseste **acquire** si **release** pe variabilele de sincronizare

- **acquire** – cand intra in sectiunea critica - datele protejate de variabila de sincronizare sunt facute consistente
- **release** - cand iese din sectiunea critica

Model specific de folosire variabile de sincronizare

- Fiecare variabila de sincronizare are un **proprietar** curent
- Proprietarul poate **modifica** datele protejate de variabila de sincronizare, in mai multe sectiuni critice (ramane proprietar)
- Un proces care vrea sa acapareze (**acquire**) trebuie sa faca o cerere proprietarului, devenind noul proprietar
- **Mod ne-exclusiv** – Mai multe procese detin variabila doar pentru **citirea** datelor protejate

Criterii de consistenta

Un proces nu poate termina un **acquire** pe o variabila de sincronizare pana cand nu s-au actualizat toate datele partajate protejate de acea variabila.

Pentru ca un proces sa modifice date partajate, el trebuie sa capate accesul exclusiv la variabila de sincronizare care protejeaza acele date (nici un alt proces sa nu detina variabila de sincronizare, nici macar in mod ne-exclusiv)

Dupa un acces exclusiv la o variabila de sincronizare, orice acces ne-exclusiv ulterior al unui alt proces la acea variabila de sincronizare este admis doar dupa ce procesul a preluat de la proprietarul variabilei de sincronizare cele mai recente copii ale variabilelor partajate protejate.

Consistenta la Intrare

P1:	Acq(Lx)	W(x)a	Acq(Ly)	W(y)b	Rel(Lx)	Rel(Ly)	
P2:					Acq(Lx)	R(x)a	R(y)NIL
P3:						Acq(Ly)	R(y)b

O secventa valida pentru consistenta la intrare.

Fiecare variabila partajata **x** este asociata cu o variabila de sincronizare (lock) **Lx**.

P2 citeste corect pe **x** (face acquire) dar nu si pe **y**

Modele de Consistenta Centrata pe Client

Este un caz special: lipsesc actualizarile simultane

Exemple:

- Sisteme de baze de date
- DNS
- WWW

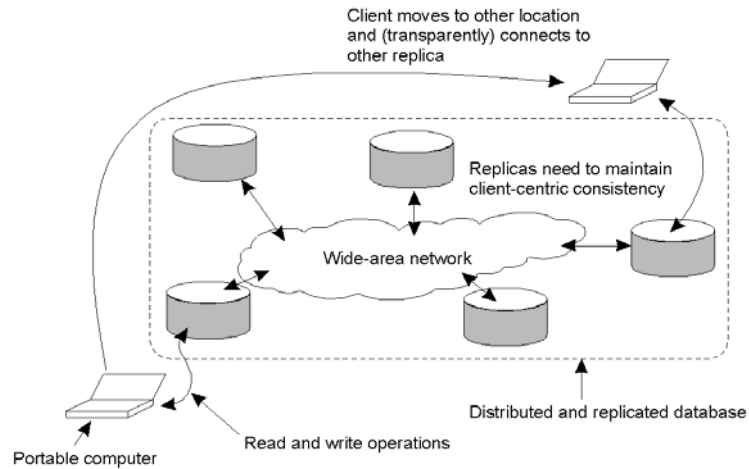
Consistenta eventuala:

- in absenta indelungata a unor actualizari, toate replicele vor deveni treptat consistente

Consistenta "Client-centric"

- un utilizator poate opera pe replice diferite
- da garantii unui singur client privind consistenta acceselor la replice diferite ale acelu client

Consistenta eventuala vs. Client-centric



Un utilizator mobil accesand diferite replici ale unei baze de date distribuite

- valorile citite si modificarile facute pe o replica sa fie regasite la cealalta replica

Citiri monotone (Monotonic Reads)

Daca un proces citeste valoarea unei date x , orice citire succesiva a lui x de catre acel proces va returna acea valoare sau o valoare mai noua.

Exemplu: baza de date distribuita, pentru e-mail.

$WS(x_1)$ este seria de operatii de scriere pe x executate la locatia 1 de la initializare; $WS(x_1; x_2) \Rightarrow WS(x_1)$ este parte a lui $WS(x_2)$

L1:	$WS(x_1)$	$R(x_1)$
L2:	$WS(x_1; x_2)$	$R(x_2)$

(a)

L1:	$WS(x_1)$	$R(x_1)$
L2:	$WS(x_2)$	$R(x_2)$ $WS(x_1; x_2)$

(b)

Operatiile *read* executate de un singur proces P pe doua copii locale diferite ale aceleiasi baze de date.

- Memorie **consistenta** "monotonic-reads"
- Memorie **ne-consistenta** "monotonic-reads"

Scriseri monotone (Monotonic Writes)

O operatie *write x* a unui proces este terminata inaintea oricarei operatii *write x* ulterioare a aceluasi proces.

Exemplu: actualizarea unei biblioteci software.

Obs. Seamana cu consistenta FIFO

L1: $W(x_1)$	L1: $W(x_1)$
<hr/>	
L2: $W(x_1)$ $W(x_2)$	L2: $W(x_2)$
(a)	(b)

Operatiile de scriere executate de un proces P pe doua copii locale diferite ale aceluasi depozit de date

- a) Un depozit consistent "monotonic-writes".
- b) Un depozit ne-consistent "monotonic-writes".

Citirea scrierilor proprii (Read Your Writes)

Efectul unei operatii *write x* a unui proces P va fi totdeauna vazut de operatiile *read x* executate ulterior de acelasi proces P .

Contra-exemplu: actualizarea paginilor Web si observarea efectelor.

Similar: actualizarea password.

L1: $W(x_1)$	L1: $W(x_1)$
<hr/>	
L2: $WS(x_1;x_2)$ $R(x_2)$	L2: $WS(x_2)$ $R(x_2)$
(a)	(b)

- a) Un depozit care **ofera** consistenta "read-your-writes".
- b) Un depozit care **nu ofera**.

Scrierile urmeaza citirile (Writes Follow Reads)

O operatie *write* x executata de un proces P dupa o operatie prealabila *read* x a aceluiasi proces se executa garantat pe aceeași valoare a lui x sau pe una mai recenta decat cea citita.

Ex.: un utilizator citeste un articol A la care raspunde cu B ; B va fi scris in orice copie a newsgroup-ului doar dupa ca A a fost scris de asemenea.

L1: $WS(x_1)$	$R(x_1)$	
<hr/>		
L2: $WS(x_1, x_2)$	$W(x_2)$	
(a)		

L1: $WS(x_1)$	$R(x_1)$	
<hr/>		
L2: $WS(x_2)$	$W(x_2)$	
(b)		

- a) Un depozit care **ofera** consistenta "writes-follow-reads".
- b) Un depozit care **nu ofera**.

Implementarea

Fiecarei operatii *write* i se asociaza un identificator global unic, de catre serverul care accepta operatia pentru prima data

Pentru fiecare client se pastreaza doua *seturi* de identificatori de *write*:

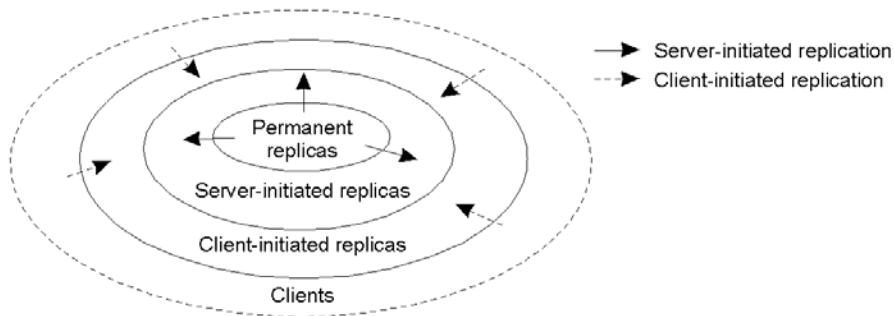
- **Read** – identificatori ai operatiilor *write* relevante pentru operatiile *read* executate de client
- **Write** – identificatori ai operatiilor *write* executate de client

Monotonic read – pentru fiecare *read*:

- Serverul primeste setul *read* al clientului
- Verifica daca toate operatiile *write* au fost facute local
- Nu -> contacteaza celelalte servere si face actualizarea
- Serverul face citirea
- Serverul actualizeaza setul *read* cu operatiile *write* realizate si care sunt relevante pentru client

Actiuni similare sunt executate pentru celelalte modele

Protocoale de distributie - Plasarea Replicilor



Organizarea logica a diferitelor categorii de copii ale depozitelor de date in trei inele concentrice.

Plasarea Replicilor (2)

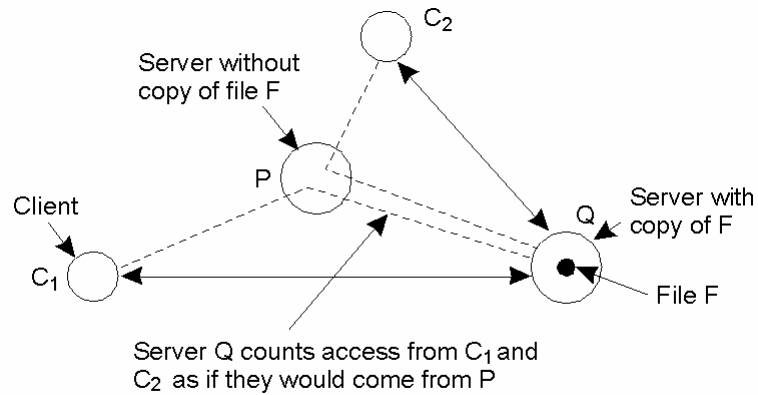
Replici permanente

1. servere "replica" situate in acelasi sistem (cluster)
2. Mirroring – replici distribuite geografic

Replici initiate de server

- Create dinamic de servere pentru a imbunatati performanta
 - replici apropiate de clienti
 - replici pentru reducerea incarcarii unui server
- Cand si unde? Algoritm (Rabinovich):
 - Fiecare server tine evidenta numarului de accese per fisier si a originii cererilor (vezi slide urmator)
 - Prag de stergere si prag de replicare
 - Intre => posibila migrare

Numararea cererilor de acces de la diversi clienti.



Plasarea Replicilor (3)

Replici initiate de clienti (caches)

- Gestionate de clienti folosind informatii de la servere
- Utilizate doar pentru a imbunatati timpul de acces la date
- Date pastrate pentru un timp limitat
- Cache-uri partajate de clienti

Propagarea Actualizarilor

Stare versus operatii

- Propagarea notificarilor
 - Utilizate de protocoalele de invalidare
 - Utilizeaza largimi de banda reduse
- Propagarea schimbarilor
 - Utila cand rata read / write este ridicata
 - Log-urile modificarilor pot fi propagate
- Propagare operatii de actualizare
 - Replicare activa
 - Cost comunicare redus

Protocoale Pull versus Push

Push:

replici permanente si "server initiated" care necesita un grad inalt de consistenta

Pull:

folosite pentru cache-uri client; utile cand rata read / write este redusa

Element	Push-based	Pull-based
Stare server	Lista replicilor si a cache-urilor client	Nimic
Message transmise	update sau invalidare plus fetch update ulterior	poll si update
Timp raspuns la client	Imediat (sau timp de fetch update)	Timp fetch update

Comparatie intre protocoale push-based si pull-based in cazul sistemelor cu mai multi clienti si un singur server.

Protocoale de actualizare hibride

Contract – o promisiune ca serverul va transmite actualizari la client pentru un timp specificat.

Timpul poate fi adaptat dinamic in functie de criteriile de **contract (lease)**

- bazat pe vechime
 - contracte de durata pentru elemente cu sansa de a ramane nemodificate
- bazat pe frecventa de innoire
 - contracte de durata pentru clienti ale caror cache-uri trebuie innoite frecvent
- overhead-ul datorat starii serverului
 - Serverul scade timpul de expirare al noilor contracte daca este supraincarcat (scade numarul de clienti)

Protocoale epidemice

Rol: propagarea actualizarilor la replici cu un numar minim de mesaje

Varianta "anti-entropie"

Categorii de servere

- **Infectios** – detine un update pe care vrea sa-l propage
- **Susceptibil** – inca ne-actualizat
- **Eliminat** – detine un update dar nu vrea sa-l propage

Abordari

- **Push based** – mai bun pentru putini infectiosi
- **Pull based** – mai bun pentru multi infectiosi

Varianta Speciala: raspandirea zvonurilor (gossiping)

- Interesul serverului de a propaga actualizarile scade pe masura ce intalneste servere actualizate deja
- Nu garanteaza actualizarea tuturor replicilor

CertIFICATE DE DECES - PROBLEME

Propagarea **stergerii** unui element de date este dificilă

Soluție: **certIFICATE DE DECES**

Când se șterge un certificat de deces?

- Execută procedura pt. a detecta dacă ștergerea este cunoscută peste tot
 - Apoi colectează certificatele de deces (similar garbage collection)
- Asociază un timp de viață maxim cu fiecare certificat
 - Există riscul de a nu atinge toate serverele
 - Soluție: un număr mic de servere păstrează nelimitat certificate "dormante"

PROTOCOALE DE CONSISTENȚĂ

Un protocol de consistență descrie o implementare a unui model de consistență

Protocoloale bazate pe o **copie primară**

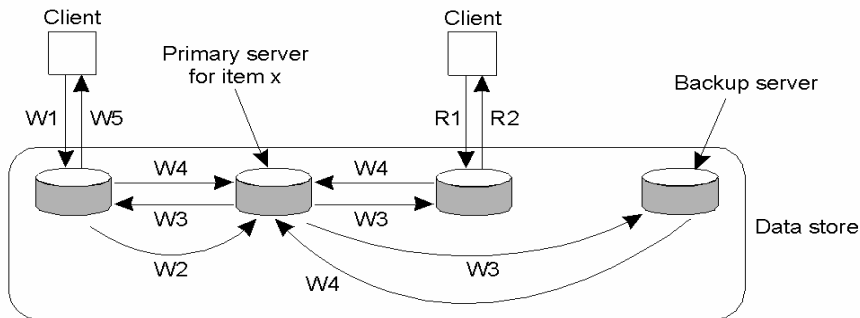
- scriere la distanță (Remote-write)
- scriere locală (Local-write)

Protocoloale cu **scriere replicată**

- replicare activă
- bazate pe quorum

Protocoloale de **coerență a cache-urilor**

Protoale Remote-Write



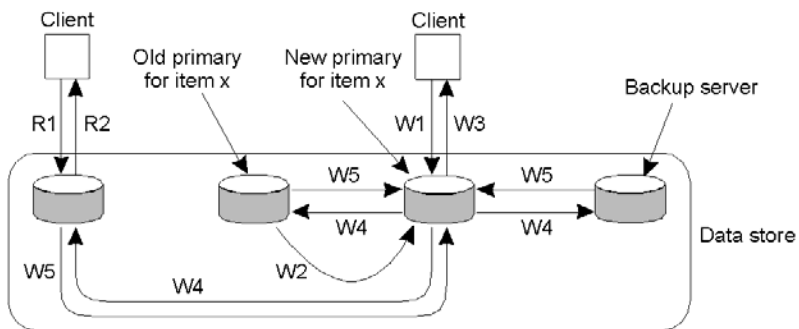
W1. Write request
 W2. Forward request to primary
 W3. Tell backups to update
 W4. Acknowledge update
 W5. Acknowledge write completed

R1. Read request
 R2. Response to read

Principiul protocolului **primary-backup**.

Implementeaza consistenta secventiala (copia primara poate ordona toate operatiile de scriere primite)

Protoale Local-Write

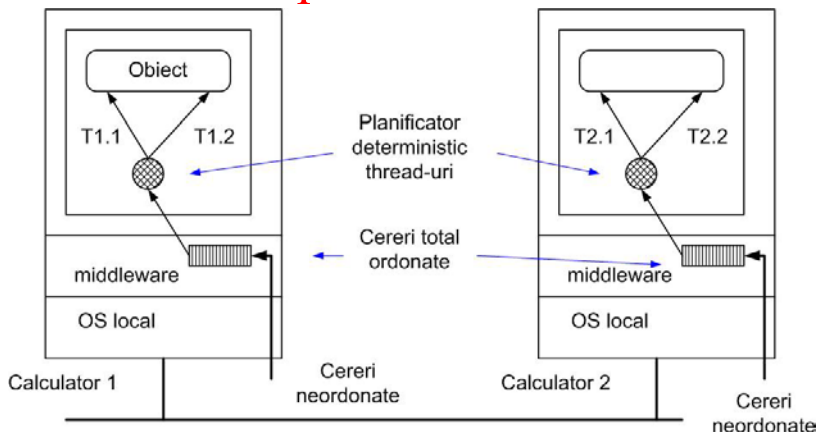


W1. Write request
 W2. Move item x to new primary
 W3. Acknowledge write completed
 W4. Tell backups to update
 W5. Acknowledge update

R1. Read request
 R2. Response to read

Protocol "**primary-backup**" in care copia primara migreaza la procesul care vrea sa faca actualizarea

Replicare activa



Doua probleme:

- Prevenirea executiei concurente a invocarilor aceluiasi obiect (**lock-uri locale** obiectului)
- Operatiile trebuie executate in aceeasi ordine in toate replicile
 - Solutia 1:
 - folosirea unei scheme **primary-based** la nivel aplicatie
 - » serializeaza cererile
 - » efort al dezvoltatorului de aplicatii

- Solutia 2: implementare in middleware

- multicast total ordonat pentru invocari
- trebuie asigurat si ca threadurile trateaza cererile in ordinea corecta

- Problema

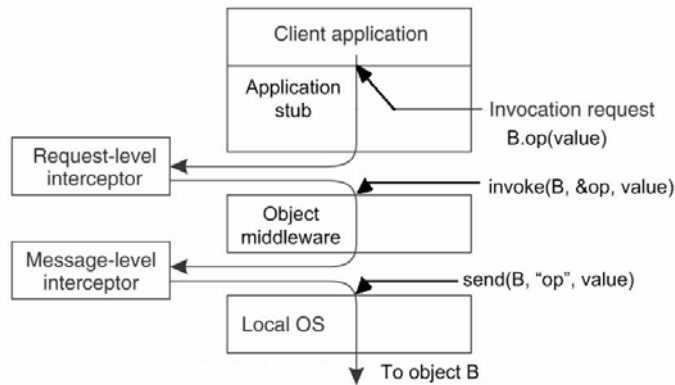
- functionare servere multithread

- preiau o cerere
- o paseaza unui thread disponibil
- trec la urmatoarea cerere
- planificatorul de threaduri aloca procesorul threadurilor executabile - **ordinea nu este aceeași in toate replicile**

- solutie: **planificator determinist**

- are la baza o varianta de **replica primara**
 - » copia primara determina ordinea threadurilor
 - » o comunica celorlalte replici
- frecvente comunicari între replici - ineficient

Separarea replicare de functionalitate obiecte

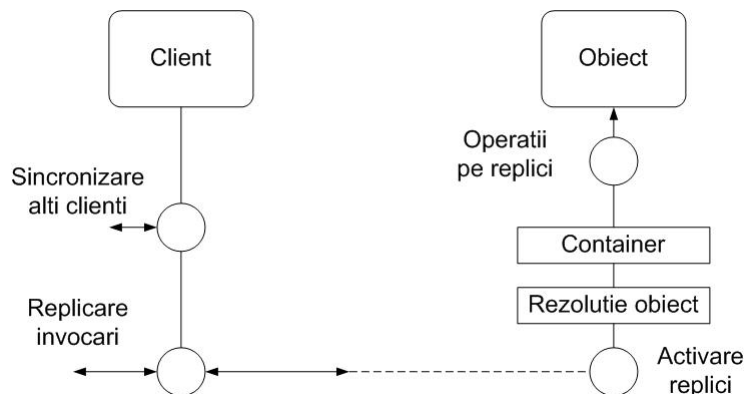


Bazata pe mecanismul [interceptorilor](#)

- constructie software care intrerupe fluxul de control normal si permite executia altor functii
- ex.: inserare interceptori in invocarea unor metode la distanta

Java IDL and Java RMI-IIOP Technologies: Using Portable Interceptors (PI)
<http://java.sun.com/j2se/1.4.2/docs/guide/idl/PI.html>

Un cadru complet de replicare



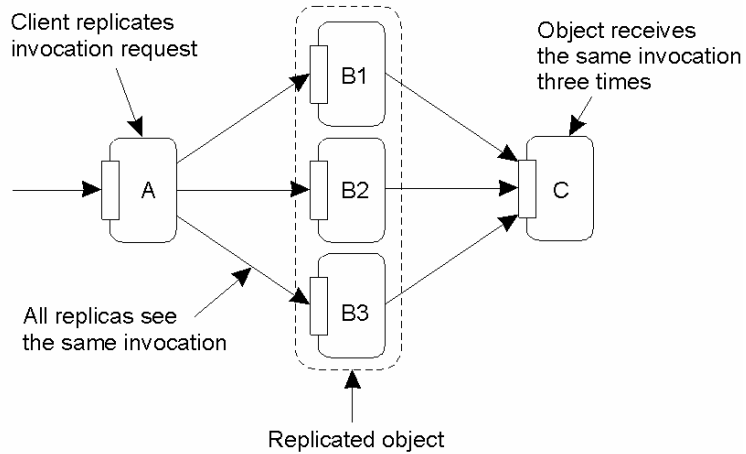
Interceptori client

- Sincronizare cu alti clienti, cand clientul este replicat
- Replicare invocari, cand invocarea trebuie transmisa mai multor replici

Interceptor server

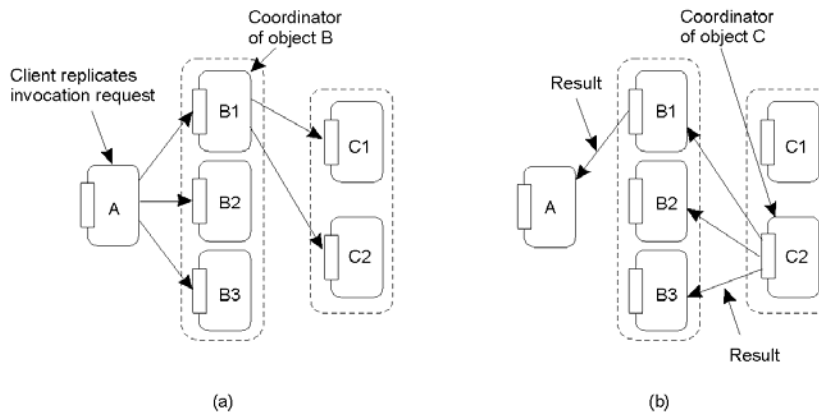
- Activare replici, daca invocarea se adreseaza mai multor replici
- Multiplicarea operatiilor daca trebuie aplicate mai multor replici

Invocari Replicate



Problema: la C ajung invocari replicate; doar una este necesara

Solutie middleware: Comunicatii constiente de replicare

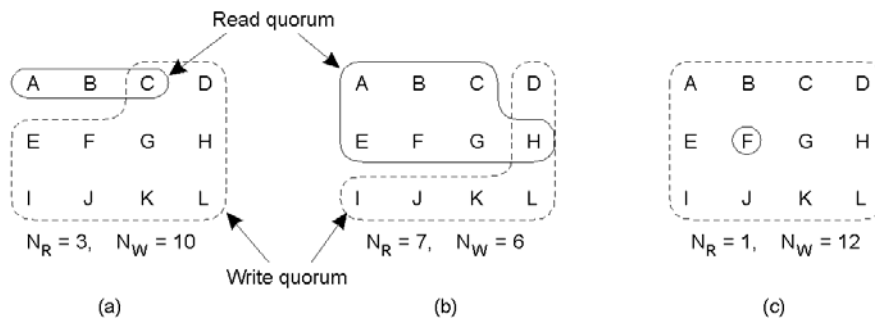


Solutie: bazata pe transmiterea multicast a invocarilor/raspunsurilor

- replicile asigneaza acelasi identificator invocarilor
- doar coordonatorul transmite invocarea
- toate replicile primesc copii ale aceluasi raspuns

Alternativa: identificarea si eliminarea invocarilor duplicate la un obiect

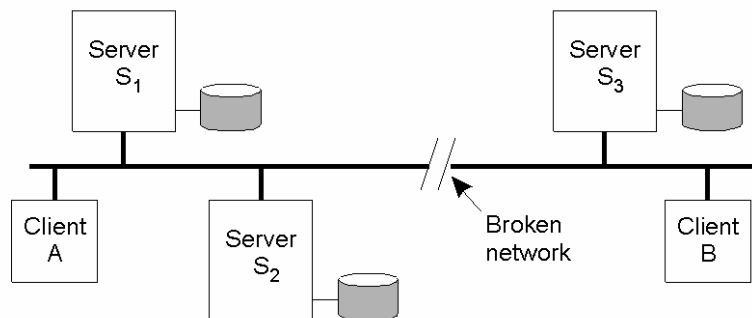
Protocoale bazate pe cvorum



Trei exemple pentru algoritmul de votare:

- Alegere corecta a seturilor de citire si scriere
- Alegere ce poate conduce la conflicte write-write
 - ex. un proces alege setul de scriere $\{A, B, C, E, F, G\}$ iar altul $\{D, H, I, J, K, L\}$
- Alegere corecta, cunoscuta ca ROWA (read one, write all)

Aplicatie: replicarea serverelor in CODA



Unitatea de replicare = volum

VSG (Volume Storage Group) = grupul de servere care au copia volumului

AVSG (Accessible Volume Storage Group) = serverele la care un client are acces; Ex.

- S_1, S_2 pentru clientul A
- S_3 pentru clientul B

CODA foloseste ROWA pentru consistenta

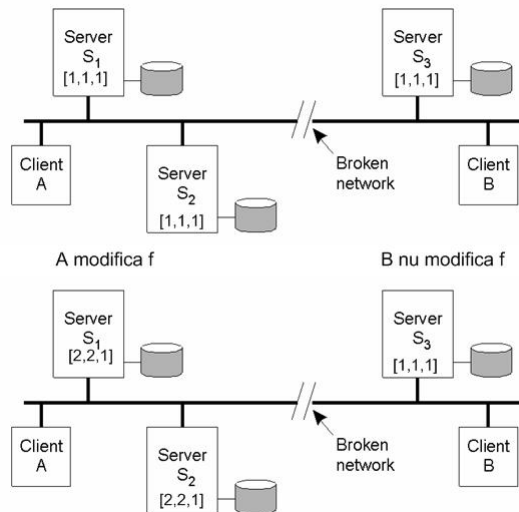
- deschide un fisier pe un membru AVSG
- la inchidere scrie toate replicile din AVSG

Varianta optimista la partitionare retea

- fiecare client lucreaza pe AVSG-ul sau

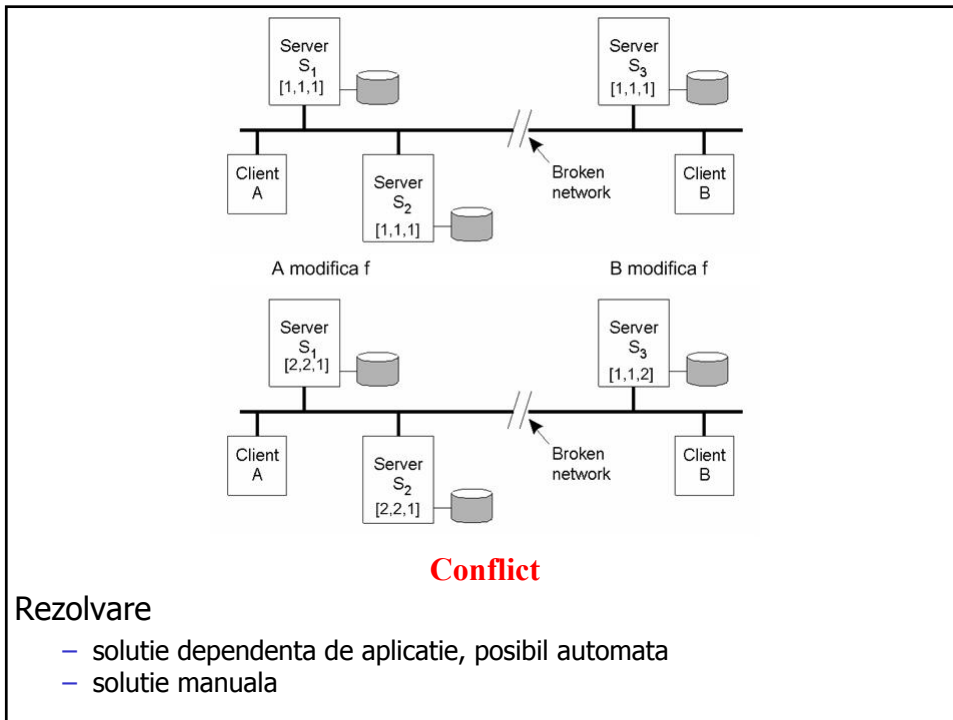
Detectie inconsistente

- CODA pastreaza pentru fiecare fisier f un vector de versiune (Coda Version Vector) $CVV_k(f)$
- $CVV_k(f)[i]$ este numarul versiunii curente a lui f pe serverul S_i , pastrat de k
- la rezolvarea defectului de partitionare a retelei, se compara CVV-urile si se detecteaza conflictele



Fara conflict

Modificarile sunt acceptate de S3



Protocoloale de coerenta cache-urilor

Caz special de replicare controlata de client

Cand sunt detectate inconsistentele? (coherence detection strategy)

- Static, in compilatoare
 - determina ce date pot conduce la inconsistente
 - insereaza cod pentru ocolirea inconsistentei
- Dinamic, la executie – trei politici
 - La accesarea unei date din cache se face mai intai verificarea
 - Lasa tranzactia sa se faca in paralel cu verificarea coerentei
 - Verifica coerenta datelor cand tranzactia se comite

Cum sunt cache-urile pastrate consistente? (coherence enforcement strategy)

- Nu permite cache-ul datelor partajate
- Lasa serverul sa transmita invalidarea sau propagarea actualizarii

Ce se intampla cand un proces modifica datele din cache?

- Serverul modifica si propaga actualizarile la cache-uri
- Clientul capata drept de acces exclusiv, modifica si paseaza actualizarile la servere (**write-through caches**)