

## **Elemente generale privind sistemul de prelucrare a tranzacțiilor.**

### **Tranzacții distribuite. Protocolul de comitere în două faze.**

(text preliminar)

Vom începe cu prezentarea unei definiții mai largi a conceptului de tranzacție, decât cea formulată mai înainte (unitate de prelucrare secvențială care respectă coerența bazei de date). Potrivit înțelegerii adoptate în (carte Gray, Reuter), o tranzacție este o colecție de operațiuni efectuate asupra stării fizice și stării abstracte a unei aplicații. Să remarcăm, în acest context, că prin aplicație se înțelege de obicei actul punerii în practică a unor principii sau idei. Definiții oarecum diferite pentru noțiunea de tranzacție corespund fie utilizatorilor finali, fie operatorilor sistemului de calcul în care sunt găzduite tranzacțiile.

Sistemul de prelucrare a tranzacțiilor (prescurtat sistemul PT) pune la dispoziție mijloacele de facilitare sau de automatizare a programării legate de aplicații, a execuției și administrării tranzacțiilor. Aplicațiile de prelucrare a tranzacțiilor evoluează în interiorul unei rețele de dispozitive, subsisteme și programe care formulează fraze de interogare și-sau operațiuni de actualizare adresate aplicației. Desigur, aplicațiile pot fi asociate unei **baze de date distribuite**. Pornind de la “intrările” menționate, aplicația întreține o bază de date ce reprezintă o stare a sistemului din lumea reală care este modelat. Firește, ne putem imagina că ieșirile furnizate de aplicație, sau reacțiile acesteia, pot conduce la schimbări ale sistemului real, respectiv a stării acestuia.

Un sistem PT conține generatoare de aplicații, rețele, baze de date, precum și aplicații. Pentru coordonarea și dirijarea fluxului de tranzacții prin sistem, este folosit un ansamblu de funcții- sau servicii-, denumit “monitor de prelucrare a tranzacțiilor”, care face parte din sistemul PT.

Este util de reamintit că **sistemele de prelucrare a tranzacțiilor au avut un rol important în evoluția calculului distribuit și a calculului tolerant la defecte**. Prin intermediul acestor sisteme a fost inițiată utilizarea informațiilor distribuite în vederea asigurării fiabilității sistemelor și aplicațiilor, a disponibilității și performanțelor acestora. Modelul client-server, foarte răspândit până nu demult, ca și apelul de la distanță al procedurilor destinat calculului distribuit, au fost de asemenea propulsate de sistemele PT. Totodată, se consideră că principalele idei unificatoare din domeniul calculului distribuit, respectiv **atomicitatea, consistența, izolarea și durabilitatea**, care împreună poartă numele de **proprietăți ACID** ale tranzacțiilor, au apărut în procesul de evoluție a

sistemelor de prelucrare a tranzacțiilor. Datorită însemnătății lor, aceste proprietăți sunt succint definite mai jos:

**Atomicitate:** schimbările stării unei tranzacții sunt atomice, în sensul că acestea au loc toate, sau nu are loc nicio schimbare. Schimbările sunt modificări ale bazei de date, mesaje și acțiuni asupra subsistemelor din mediul extern.

**Consistență:** o tranzacție constituie o transformare corectă a stării bazei de date. Acțiunile întreprinse în grup nu încalcă niciuna din restricțiile de integritate specifice stării menționate. De aici rezultă că tranzacției trebuie să îi corespundă un program corect.

**Izolare:** chiar dacă tranzacțiile sunt executate în mod concurrent, fiecare tranzacție T “vede” execuția celorlalte tranzacții ca fiind produsă fie anterior față de T, fie după T, dar nu ambele.

**Durabilitate:** odată ce o tranzacție a fost executată cu succes (deci a fost “completată” sau “a comis”), schimbările pe care ea le-a produs în starea bazei de date nu vor fi afectate de defecțiuni.

Evoluția globală a unei tranzacții, cu respectarea proprietăților ACID, poate fi descrisă după cum urmează. Programul de aplicație declară debutul unei tranzacții noi prin comanda Begin-Work(). În continuare, toate operațiile efectuate de program vor face parte din această tranzacție. În același timp, toate operațiunile executate de alte programe în legătură cu programul de aplicație sunt parte a tranzacției de bază. Prin inițierea comenzii Commit\_Work, programul declară că tranzacția constituie o transformare completă și corectă. Rezultatele execuției tranzacției – sau efectele acesteia- devin durabile, de îndată ce tranzacția comite.

Pe de altă parte, este posibil ca pe parcursul execuției tranzacției unele operațiuni să fie eronate. În acest caz, aplicația poate „distruge” toate operațiunile, prin comanda Rollback\_Work(). Mai observăm că atunci când în timpul execuției unei tranzacții are loc o defecțiune de natură tehnică, sistemul de gestiune poate ”trimite înapoi”, sau „returna” tranzacția. Pentru încadrarea transformărilor ACID se folosesc comenzile Begin-Commit sau Begin-Rollback.

Elementele de mai sus reprezintă o abordare convenabilă pentru realizarea **aplicațiilor distribuite**. Fiecare modul al aplicației va fi o tranzacție sau o sub-tranzacție. Atunci când execuția se desfășoară cu succes, tranzacția comite și toate modulele trec în o stare durabilă nouă.. În situația în care apar erori sau alte incidente, modulele tranzacției vor fi readuse automat în starea inițială, cea de la începutul execuției tranzacției. Întrucât operațiunile de comitere și de revenire se bazează pe o logică „automată”, realizarea modulelor reclamă folosirea de metode cu semantică elementară pentru descrierea defecțiunilor.

## Tranzacții distribuite

În cadrul unei **baze de date distribuite în nodurile unei rețele**, accesul la multitudinea de obiecte de date se realizează de obicei prin intermediul tranzacțiilor. Se înțelege că tranzacțiile trebuie să respecte regulile ACID.

Deoarece **tranzacțiile care sunt executate într-un sistem distribuit pot prelucra fie date rezidente într-un singur loc, fie date amplasate în sisteme distincte**, vom distinge **tranzacții locale** și **tranzacții globale**. Primele au acces la informațiile conținute în o singură bază de date locală și pot prelucra aceste informații. Tranzacțiile globale realizează accesul la informații amplasate în mai multe baze de date locale și, desigur, pot prelucra aceste informații.

Respectarea proprietăților ACID în procesul execuției **tranzacțiilor locale** reclamă, practic, participarea tuturor componentelor sistemului de gestiune a bazei de date implicate în gestiunea tranzacțiilor. Asigurarea consistenței unei tranzacții oarecare trebuie dată de programatorul de aplicație, care codifică tranzacția. Pentru respectarea atomicității tranzacției, chiar dacă apar defecțiuni, se poate folosi tehnica de modificare „amânată”, sau tehnica modificării imediate a bazei de date. În primul caz, modificările aduse bazei de date sunt înregistrate într-un „jurnal”, însă execuția tuturor operațiilor de înregistrare din tranzacție este amânată până la comiterea parțială a tranzacției (reamintim că prin comitere parțială înțelegem situația în care acțiunea finală a tranzacției a fost executată). Modificarea imediată a bazei de date constă în efectuarea schimbărilor bazei de date atunci când tranzacția este încă în stare activă; schimbările aduse de tranzacții ce se găsesc în stare activă se numesc **modificări ne-comise**. Responsabilitatea operațiilor de mai sus revine **controlorului** (sau „manager”-ului) **de tranzacții**. În fapt, acest controlor, prin observarea memoriei secundare (disc), urmărește valorile vechi ale tuturor obiectelor de date asupra cărora tranzacția a efectuat înregistrări, iar dacă tranzacția nu a completat execuția sa, restaurează vechile valori astfel încât aparent tranzacția nu ar fi fost niciodată executată. Durabilitatea este în sarcina **controlorului de revenire** (sau de recuperare). Prin grja acestui sub-sistem, modificările aduse bazei de date de către o tranzacție sunt înregistrate pe disc înaintea completării tranzacției. De asemenea, este asigurată disponibilitatea unei informații suficiente privind modificările efectuate și înregistrate pe disc, astfel încât actualizările făcute să poată fi reconstruite atunci când sistemul de baze de date este repornit după o defecțiune. În sfârșit, **controlorul concurenței** este sub-sistemul care, printre altele, asigură respectarea proprietății de izolare. Printre tipurile de planificare ce contribuie la asigurarea acestei proprietăți menționăm serializabilitatea de conflict și serializabilitate de vedere.

În cazul **tranzacțiilor globale**, menținerea proprietăților ACID este sensibil mai complicată, având în vedere faptul că la execuția unei planificări de tranzacții iau parte mai multe noduri ale unei

rețele. Defecțiunile apărute în unul sau mai multe din noduri, ca și defecțiunile canalelor de comunicație, vor conduce aproape sigur la prelucrări incorecte ale informațiilor.

Din arhitectura unei **baze de date distribuite** fac parte două componente principale: **controlorul de tranzacții și coordonatorul tranzacțiilor**. Aceste sub-sisteme sunt **amplasate în nodurile rețelei de calculatoare**.

Funcțiunea **controlorului de tranzacții** constă în **respectarea proprietăților ACID** ale tranzacțiilor care prelucrează informații memorate **în nodul cu care controlorul este asociat**. Deoarece avem în vedere o bază de date distribuită, se înțelege că în rețea vor fi executate și tranzacții globale, care efectuează operațiuni asupra datelor depuse în mai multe noduri. În acest proces vor colabora controloare de tranzacții amplasate în nodurile corespunzătoare. Într-un nod dat al rețelei, în care se desfășoară prelucrarea concurentă a unui număr de tranzacții, **controlorul de tranzacții ocupă un loc specific în activitatea algoritmului de control al concurenței**. De asemenea, acest sub-sistem are **responsabilitatea conținutului jurnalului** menținut în vederea revenirii după defecțiuni.

La rândul său, **coordonatorul de tranzacții asociat unui nod, dirijează execuția tranzacțiilor locale și globale ce sunt inițiate și funcționează în nodul respectiv**. **Funcția de coordonare menționată este specifică numai mediului distribuit**. În sarcina coordonatorului intră **inițierea execuției fiecărei tranzacții, segmentarea tranzacțiilor** în scopul distribuirii sub-tranzacțiilor obținute în noduri adecvate, precum și **coordonarea încheierii activității fiecărei tranzacții**. Subliniem faptul că tranzacțiile globale fie comit în toate nodurile în care sunt executate sub-tranzacțiile din care sunt formate, fie sunt abandonate în toate nodurile.

În legătură cu defectele ce apar uneori în sistemele centralizate, inclusiv în cele care prelucrează baze de date, menționăm erorile de echipament („hardware”), erorile programelor, deteriorarea sistemelor de discuri magnetice. Tema defecțiunilor posibile într-un **sistem distribuit** care găzduște o bază de date are particularități evidente. Deteriorarea componentelor unui nod, defecțiuni ale liniilor de comunicație, perturbarea manierei inițiale de partiționare a rețelei, coruperea mesajelor prin erori sunt tipurile principale de defecte ce au loc în mediu distribuit. Atunci când se defectează mediul de comunicație, mesajele transmise între noduri sunt de obicei re-rutate. În ipoteza dispariției conexiunii între două noduri, rețeaua va fi partiționată din nou. Erorile care apar în mesaje sunt tratate de protocoalele de control al procesului de comunicație, de exemplu TCP/IP.

Într-un paragraf imediat precedent, am afirmat că la sfârșitul execuției unei tranzacții globale- sau distribuite-, aceasta trebuie fie să comită, fie să fie abandonată în toate nodurile rețelei. În o

formulare mai puțin tehnică, rezultatul final al execuției tranzacției trebuie să întrunească „acordul” tuturor nodurilor, fapt în care constă respectarea proprietății de atomicitate a tranzacției.

Pentru asigurarea proprietății de atomicitate, coordonatorul tranzacției menționate este prevăzut cu un **protocol de comitere**. Unul dintre sistemele de programe cele mai frecvent folosite, din această clasă, este **protocolul de comitere în două faze (prescurtat, 2PC)**. O variantă mai complexă este reprezentată de protocolul de comitere în trei faze (3PC).

### **Protocolul de comitere în două faze (2P Commit) - structură, funcționare.**

În cele ce urmează, **în contextul unei baze de date distribuite**, va fi descrisă activitatea protocolului 2PC până la finalizarea unei tranzacții, respectiv până la comiterea acesteia. De asemenea, protocolul 2PC îndeplinește funcții specifice în situațiile de defect, de revenire după defecte și de control al concurenței tranzacțiilor. Modul de realizare a acestor ultime funcții va fi prezentat ulterior.

#### **Acțiunea de comitere.**

Fie o **tranzacție T** care intră în lucru în **nodul Ni**, unde **coordonator este TCi**. Vom distinge două faze în desfășurarea procesului de încheiere a execuției tranzacției T. Așa cum se va vedea, procesul se poate sfârși fie prin comiterea tranzacției, fie prin abandonarea acesteia. Protocolul 2PC intră în lucru atunci când, prin decizia coordonatorului TCi, tranzacția T încheie (sau completează) execuția sa, respectiv atunci când toate nodurile în care a fost prelucrată T anunță coordonatorul TCi că tranzacția a fost completată.

**Prima fază.** Odată cu intarea în faza de început a procesului, coordonatorul TCi introduce în jurnal înregistrarea cu conținutul <prepare T> și, totodată, înmagazinează (sau, cu un termen folosit frecvent, „forțează”) conținutul jurnalului în memoria „stabilă” (pe discuri magnetice) a sistemului. În continuare, TCi transmite mesajul privind pregătirea tranzacției tuturor nodurilor în care aceasta a fost prelucrată. În aceste noduri, în urma recepționării mesajului, fiecare controlor de tranzacții asociat decide asupra comiterii porțiunii de tranzacție care îi corespunde. În situația în care răspunsul este afirmativ, controlorul de tranzacții din nod introduce în jurnal înregistrarea cu conținutul <ready T>, iar apoi „forțează” conținutul jurnalului în memoria stabilă, inclusiv toate înregistrările din jurnal care se referă la tranzacția T. Apoi, controlorul de tranzacții transmite coordonatorului TCi mesajul „ready T”. Atunci când decizia este negativă, controlorul de tranzacții din nod înscrie înregistrarea <no T > în jurnal, în continuare transmițând coordonatorului TCi mesajul „abort T”.

**A doua fază.** Coordonatorul de tranzacții TCi poate stabili dacă tranzacția T va fi comisă sau abandonată, de îndată ce a primit răspunsuri de la toate nodurile la mesajul său de pregătire a

tranzacției. Decizia respectivă poate fi luată de TCi și atunci când de la transmiterea mesajului <prepare T> a trecut un timp dat.

Dacă toate nodurile ce participă în proces au răspund prin mesajul „ready T”, tranzacția T poate fi comisă, în jurnal se introduce înregistrarea <commit T>, iar conținutul jurnalului este înmagazinat în memoria stabilă. În situația contrară, T este abandonată, urmând ca în jurnal să fie înscrisă informația <abort T>, iar jurnalul- forțat în memoria stabilă. Potrivit exprimării din literatura de specialitate (vezi Silberschatz, pg 216), în acest moment evoluția tranzacției este definitivă. Ca urmare, sistemul coordonator TCi transmite tuturor nodurilor ce au luat parte la execuție mesajul „commit T” sau „abort T”. Primind unul din aceste mesaje, nodurile îl înscriu în jurnal.

Să adăugăm câteva comentarii față de elementele de mai sus. În primul rând, de îndată ce controlorul de tranzacții emite un mesaj „ready T” către coordonator, tranzacția se va găsi în situația numită „ready state” în nodul corespunzător. În fapt, mesajul transmis coordonatorului este o „promisiune” formulată la nivelul nodului privind invitația de comitere sau de abandonare exprimată de TCi. În acest context, trebuie remarcat că pentru a formula un asemenea răspuns, este necesară înscrierea unei informații corespunzătoare în memoria stabilă. Dacă informația menționată nu va fi disponibilă, este posibil ca în situația în care, în urma apariției unui defect, nodul în cauză iese din funcțiune după ce a transmis mesajul „ready T”, intenția de comitere sau abandonare să nu mai poată fi realizată în mod corespunzător.

În al doilea rând, tot aici vom mai observa că drepturile de închidere ale tranzacției rămân valabile până la comiterea acesteia, precum și faptul că T poate fi abandonată în orice moment, anterior celui în care controlorul de tranzacții din nod trimite coordonatorului mesajul „ready T”.

În al treilea rând, vom sublinia faptul că de îndată ce un nod a declarat abandonarea tranzacției T prin răspunsul „abort T”, evoluția tranzacției este definitivă, drept urmare a condiției de unanimitate pentru comitere. În legătură cu aceasta, observăm că decizia de abandonare a tranzacției T poate fi luată în mod unilateral de către coordonatorul TCi, al nodului în care s-a produs prelucrarea tranzacției. Decizia finală asupra evoluției tranzacției T este luată în momentul în care TCi înscrie hotărârea de comitere sau abandonare adoptată în jurnal, apoi în memoria stabilă. Unele implementări ale protocolului 2PC prevăd ca, la sfârșitul celei de a doua faze, controlorul de tranzacții din nod să transmită coordonatorului mesajul „acknowledge T”, urmând ca după recepționarea mesajului din partea tuturor nodurilor, TCi să introducă în jurnal înregistrarea <complete T>.

Funcționarea protocolului de comitere în două faze poate fi descrisă în o manieră alternativă, prezentată succint în continuare. Fie o bază de date distribuită în N noduri. Mai jos, vor fi folosite notațiile: <.....> pentru reprezentarea înregistrărilor, iar “.....” pentru mesaje. Admitem că din

componenta unui nod “i” fac parte o unitate centrală de prelucrare UC<sub>Pi</sub>, o memorie centrală MC<sub>i</sub>, o memorie de tip cache MCH<sub>i</sub>, în nod fiind înmagazinat codul tranzacțiilor locale și subtranzacțiilor ce aparțin unor tranzacții globale, precum și sistemele de programe corespunzătoare controlorului Tc<sub>tri</sub>, respectiv coordonatorului de tranzacții TC<sub>i</sub>. La elementele de mai sus se adaugă memoria secundară MSEC<sub>i</sub>. În continuare, este prezentată schematic o desfășurare posibilă a activității protocolului 2P Commit, la momentele de timp t<sub>1</sub>, t<sub>2</sub>, etc.

### Prima fază

- t<sub>1</sub> TC<sub>i</sub> <prepareT> -> jurnal
- t<sub>2</sub> TC<sub>i</sub> <jurnal> -> memoria stabilă MSEC<sub>i</sub>
- t<sub>3</sub> TC<sub>i</sub> <prepareT> -> noduri
- t<sub>4</sub> Tc<sub>trj</sub> decizie privind comiterea porțiunii proprii de tranzacție
- t<sub>5</sub> decizie = Da, Tc<sub>trj</sub> <readyT> -> memoria stabilă,  
sau  
decizie = Nu, Tc<sub>trj</sub> <noT> -> jurnal
- t<sub>6</sub> decizie = Da, Tc<sub>trj</sub> mesaj “readyT” -> TC<sub>i</sub>,  
sau  
decizie = Nu, Tc<sub>trj</sub> mesaj “abort” -> TC<sub>i</sub>.

### A doua fază

- t<sub>7</sub> TC<sub>i</sub> analizează răspunsurile Tc<sub>trj</sub> din momentul t<sub>6</sub>
- t<sub>8</sub> toate Tc<sub>trj</sub> -> mesaj „readyT”, TC<sub>i</sub> <commitT> -> jurnal,  
sau  
nu toate Tc<sub>trj</sub> -> mesaj „ready”, TC<sub>i</sub> <abort> -> jurnal
- t<sub>9</sub> TC<sub>i</sub> <jurnal> -> memoria stabilă MSEC<sub>j</sub>
- t<sub>10</sub> TC<sub>i</sub> <commitT> în toate nodurile, sau <abortT> în toate nodurile
- t<sub>11</sub> fiecare Tc<sub>trj</sub> participant în prelucrarea tranzacției <commitT> -> jurnal sau  
<abortT> -> jurnal.

În ipoteza că unul sau mai multe din nodurile rețelei (respectiv dintre controloarele de tranzacție) nu „votează” favorabil, tranzacția nu poate comite. În această situație, pentru tranzacție trebuie să se ajungă în final la decizia „Abort\_Work()”. Jurnalul tranzacției va fi citit înapoi, sub controlul coordonatorului de tranzacții. În acest proces, pentru fiecare înregistrare din jurnaș coordonatorul apelează nodul (controlorul) care a operat ultima actualizare, determinându-l să inițieze o operațiune de „distrugere” a informației respective. Atunci când în procesul inspecției „inverse” a jurnalului se ajunge în punctul de început al tranzacției, respectiv acolo unde trebuie să se afle comanda Begin\_Work(), va fi efectuată operațiunea de „distrugere” a acestei ultime comenzi, apoi în jurnal se introduce o înregistrare de abandonare. În sfârșit, în jurnal este introdusă o nouă înregistrare, care consemnează „completarea” abandonării. Mai menționăm că jurnalului tranzacției i se asociază un controlor propriu, a cărui sarcină este de a gestiona înregistrările din jurnal, prin intermediul unei liste înlănțuite, în care fiecărei înregistrări i se atribue o cheie unică. Unul din obiectivele pentru care în procesul de inspecție a conținutului jurnalului se folosește o listă înlănțuită este asigurarea eficienței operațiunilor de „distrugere” evocate mai sus.