

Optimizing Queries in Distributed Systems

Prof.dr. Ion LUNGU, Anca Georgiana FODOR
Catedra de Informatică Economică, A.S.E. București

This research presents the main elements of query optimizations in distributed systems. First, data architecture according with system level architecture in a distributed environment is presented. Then the architecture of a distributed database management system (DDBMS) is described on conceptual level followed by the presentation of the distributed query execution steps on these information systems. The research ends with presentation of some aspects of distributed database query optimization and strategies used for that.

Keywords: distributed system, distributed system architecture, distributed query, distributed query optimization.

Sisteme de baze de date distribuite

Apariția sistemelor distribuite a fost determinată de modul de desfășurare a activității companiilor cu sedii aflate la distanțe considerabile unele de altele, uneori chiar pe continente diferite, structura organizațională promovând un mod de lucru descentralizat. Sistemele distribuite ajută la rezolvarea problemei “insulelor informaționale”, apărute ca urmare a separării geografice dar și datorită diversității arhitecturilor calculatoarelor și a protocoalelor de comunicație.

În prezent, sistemele de baze de date distribuite sunt organizate ca sisteme software multifir, conform [SMBL02]. Arhitectura generică pentru sistemele de baze de date distribuite este arhitectura pe trei niveluri fiind alcătuită din componenta de gestiune a bazei de date, componenta de gestiune a proceselor și aplicațiilor și respectiv componenta de interfață cu utilizatorii. Similar acestei arhitecturi software cu trei componente, există **arhitectura datelor**, pe trei niveluri, prezentată în figura 1.

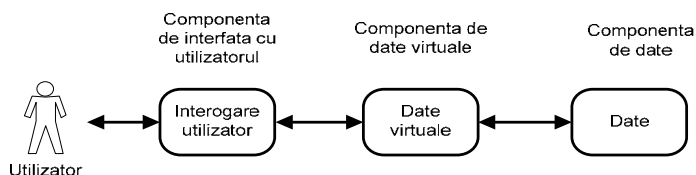


Figura 1. Arhitectura datelor pe trei niveluri

Componenta de date este alcătuită din datele stocate în sistemul distribuit care sunt disponibile în infrastructura rețelei de calculatoare. *Componenta de date virtuale* asigură transparența distribuției datelor, oferind utilizatorului accesul unitar la date. *Componenta interfeței cu utilizatorul* permite acestuia să formuleze interogări asupra uneia sau mai multor view-uri.

Responsabilitățile funcționale din arhitectura datelor prezentată mai sus revin unor sisteme de calcul, așa cum se vede în figura 2. Componenta de date este reprezentată de unul sau mai multe sisteme de gestiune a bazelor de date (SGBD) care funcționează pe mai multe servere de baze de date. Datele virtuale sunt

gestionate de componente software ce rulează pe un *server de aplicații*. *Calculatorul client* este punctul de acces al utilizatorului la sistem și modelează componenta de interfață cu utilizatorul din arhitectura datelor. Funcția acestuia este doar de a prelua interogările utilizatorului și a-i prezenta acestuia rezultatul cererilor sale. În unele cazuri, părți din procesarea interogărilor pot fi direcționate către calculatorul client, dar în general, acesta are mai mult rolul de comunicare decât de prelucrare.

Toate cele trei componente manipulează date. Fiecare componentă accesează datele de la componenta anterioară și le transferă componentei următoare. În final utilizatorul va

primi datele prin intermediul componentei "Interogare utilizator" care funcționează pe calculatorul client.

În cazul sistemelor în care componenta de

date virtuale lipsește, distribuția datelor nu este ascunsă utilizatorului, și acesta trebuie să formuleze interogări direct SGBD-urilor gestionându-și singur distribuția datelor.

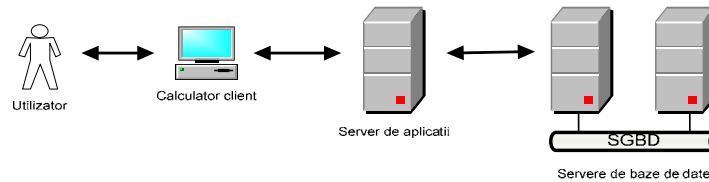


Figura 2. Arhitectura de sistem pe trei niveluri

Arhitectura de referință a unui SGBDD

Sistemele de gestiune a bazelor de date distribuite (SGBDD) promovează o arhitectură multinivel pentru a obține transparența distribuției, așa cum se poate vedea și în figura 3. Arhitectura SGBDD de tip ANSI/SPARC es-

te o arhitectură de referință pentru bazele de date distribuite și prezintă conceptual o bază de date distribuită (BDD). Această arhitectură nu este implementată explicit în toate sistemele de baze de date distribuite.

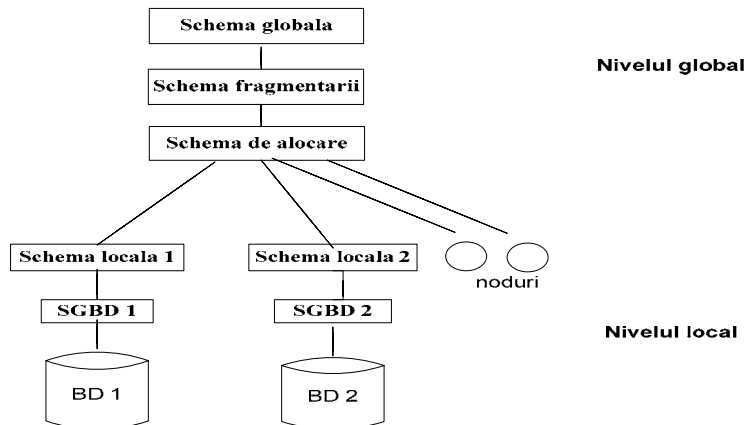


Figura 3 Arhitectura de referință a SGBDD

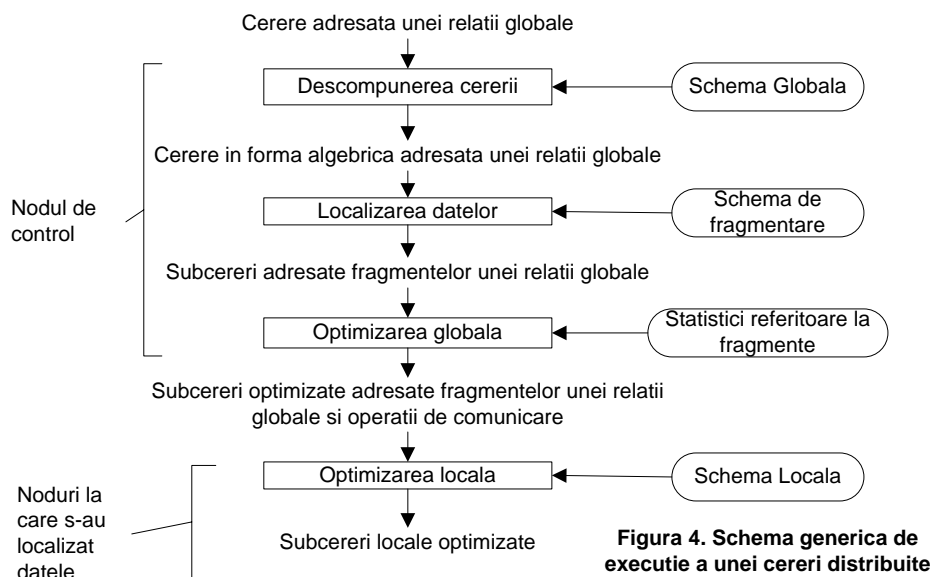
Într-un SGBDD, datele sunt organizate pe două nivele: nivelul global și nivelul local. **Nivelul global** permite integrarea bazelor de date locale într-o bază de date globală prin intermediul schemei globale, schemei de fragmentare și schemei de alocare. *Schema globală* definește toate datele conținute în BDD și este descrisă de un set de relații globale. Fiecare relație globală poate fi împărțită în mai multe părți disjuncte numite **fragmente**. Corespondența între relația globală și fragmente este definită în *schema de fragmentare*. Fragmentele sunt porțiuni logice ale relațiilor globale care pot fi alocate fizic pe una sau mai multe stații ale rețelei. *Schema de alocare* definește nodurile unde sunt alocate fragmentele. **Nivelul local** permite tratarea fiecărei baze de date locale ca o bază de date centralizată. La acest nivel *schema loca-*

lă – care depinde de tipul SGBD local - realizează corespondența între imaginea fizică a relațiilor globale la nod și obiectele manipulate de SGBD-ul local.

Execuția cererilor în sistemele de baze de date distribuite

Conform [OSVA95], în sistemele distribuite, *obiectivul execuției cererilor* este de a transforma o cerere globală, adresată bazei de date văzută ca un întreg de către utilizatorii ei, în comenzi de manipulare a datelor de nivel scăzut, adresate bazelor de date locale, aflate la noduri, folosind o strategie de execuție eficientă.

Execuția unei cereri distribuite este compusă din mai multe activități, corespunzătoare etapelor generice ale procesării cererii. Conform [OSVA95], aceste etape se desfășoară ca în figura 4.



Execuția unei cereri pornește de la cererea globală exprimată folosind operatorii de calcul relațional. Această cerere este adresată unor relații globale, distribuția datelor fiind invizibilă la acest nivel. Ea este descompusă în subcereri exprimate cu ajutorul operatorilor algebrei relaționale. Aceste subcereri se transformă din subcereri globale în cereri adresate fragmentelor relațiilor globale, folosind informații despre schema de fragmentare. Optimizarea globală constă în găsirea celei mai bune ordonări a operațiilor în cadrul subcererilor adresate fragmentelor, incluzând și operațiunile de comunicație care minimizează o funcție de cost. Optimizarea locală a fiecărei subcereri la un anumit nod se realizează folosind informații din schema locală. Aceste patru etape realizează descompunerea cererii globale într-o secvență de operațiuni locale optimizate, fiecare din acestea acționând asupra unei baze de date locale. Primele trei se desfășoară la un nod de control, care deține informații despre schema globală, ultima etapă fiind executată la nodurile unde sunt localizate datele necesare procesării cererii.

Factori care influențează optimizarea cererilor distribuite

O etapă importantă în execuția cererilor distribuite o reprezintă optimizarea execuției acestora. Deși există mai multe variante de transformare a aceleiași cereri globale, este esențial să reținem ca soluție cea variantă care optimizează (minimizează) consumul de

resurse cu prețul unui efort rezonabil.

Interogarea rezultată în urma etapei de localizare a datelor este o interogare adresată fragmentelor exprimată în algebra relațională și are o formă ce poate fi executată, elementele neesențiale fiind deja eliminate.

Este totuși necesară o optimizare a acestei forme a cererii ținând cont de dimensiunile fragmentelor implicate în execuție, de caracteristicile rețelei de comunicație peste care funcționează sistemul distribuit, precum și de resursele fizice ale acestuia (capacitate de procesare, memorie disponibilă și spațiu pe discuri).

Optimizarea poate fi privită ca activitatea prin care se încearcă obținerea celei mai bune soluții posibile în condițiile de mediu date. Spun posibil, deoarece nu întotdeauna se poate obține soluția “cea mai bună”, în valoare absolută, din motive lesne de înțeles, și anume, sau ar necesita un timp prea mare de căutare în spațiul soluțiilor, sau ar bloca prea multe resurse, comparativ cu o soluție apropiată de “cea mai bună”, dar care s-ar putea obține cu eforturi vizibil mai mici.

Problema optimizării execuției unei interogări distribuite are ca *obiectiv major* obținerea unui timp de răspuns al sistemului cât mai mic în condițiile minimizării costului total de execuție.

Costul total de execuție al unei interogări într-un sistem distribuit, conform [OSVA95], se compune din costul de procesare, costul de accesare a discurilor și costul comunicației

între noduri. Costul comunicației se compune din costul trimitere și primire mesaje între noduri și respectiv costul de transfer date între noduri. Costurile de comunicație depind în mare măsură de tipul rețelei de comunicație. În cazul unei rețele LAN, putem considera costul unitar de transfer de date identic între oricare două noduri. În cazul unei rețele WAN, costurile unitare de transfer date între noduri sunt diferite, în funcție de distanța între noduri care poate varia semnificativ de la o pereche de noduri la alta. De asemenea, în cazul WAN, costurile de procesare și respectiv costurile de accesare a discurilor sunt nesemnificative față de cele de comunicație și deci se pot neglija, având funcția de cost total compusă doar din costul de transfer date și costul de transmitere/primire de mesaje. Formula costului total se poate scrie astfel [OSVA95]:

$$\text{Cost total} = C_{CPU} * nr_instrucțiuni + C_{I/O} * nr_I/O + C_{MSG} * nr_mesaje + C_{TR} * cantitate_date$$

Timpul de răspuns al sistemului distribuit se calculează din momentul lansării în execuție a interogării până la primirea răspunsului din partea sistemului. O influență semnificativă în minimizarea timpului de răspuns o au procesarea paralelă și respectiv comunicația paralelă. Cu cât crește execuția și respectiv comunicarea în paralel, cu atât timpul de răspuns va fi mai mic. Și în cazul timpului de răspuns este foarte important tipul rețelei WAN sau LAN, deoarece, timpul de execuție a cererii, precum și timpul necesar accesului la disc devin nesemnificative față de timpul necesar transferului de date și cel al transmiterii de mesaje, în cazul rețelelor WAN, spre deosebire de rețelele LAN. Formula timpului de răspuns al sistemului distribuit se scrie, conform [OSVA95]:

$$\text{Timp răsp} = T_{CPU} * nr_instrucțiuni + T_{I/O} * nr_I/O + T_{MSG} * nr_mesaje + T_{TR} * cantitate_date$$

După cum am precizat, minimizarea timpului de răspuns se realizează prin creșterea gradului de paralelism al execuției unei interogări distribuite. Aceasta nu implică neapărat că și costul total va fi minimizat. Din contră, costul total poate crește încercând să sporim

gradul de paralelism al execuției și transmisiiei. Pe de altă parte, minimizarea costului total implică îmbunătățirea utilizării resurselor, în detrimentul timpului de răspuns, care se va mări. În practică este de dorit un compromis între cele două obiective.

Într-unul din studiile de specialitate [KAMO94] se descriu două *strategii de optimizare a execuției interogărilor* atunci când sunt necesare date de la mai multe noduri. Acestea strategii sunt: *query-site* și *move-small*.

Strategia *query-site* spune că dacă avem o interogare care folosește date ce se regăsesc în mai multe fragmente la noduri diferite, cel mai "ieftin" ar fi să transferăm toate acele fragmente la nodul de unde s-a inițiat interogarea și să o executăm acolo.

Strategia *move-small* spune că dacă o operație binară implică două fragmente ce se află la noduri diferite, atunci transferăm fragmentul mai mic la nodul unde este localizat fragmentul mai mare.

Este vorba desigur de rezultatele intermediare care s-au obținut în urma operațiilor unare (selecție, proiecție) și care urmează să participe la operațiile de "join" sau reuniune, mari consumatoare de resurse. Este important să ținem cont de aceste strategii de execuție pentru a minimiza cantitatea de date transferate între noduri, și implicit, costul de comunicație, care este o parte semnificativă a costului total, în special în rețele de tip WAN.

Componentele optimizării cererilor distribuite, conform [MULL96], sunt metoda de acces la date, criteriile de uniune folosite și costurile de comunicație.

Majoritatea SGBD-urilor permit accesarea tabelor în două moduri: secvențial sau indexat. Cea mai potrivită metodă de acces la date poate fi determinată de contextul cererii. Este posibil ca aceasta să returneze 90% din tuplurile tabelii, caz în care vom folosi accesul secvențial. În cazul în care cererea va returna 10% din tupluri, folosirea unui index ar crește viteza de răspuns a sistemului.

Dacă interogarea accesează mai multe tabele cărora li se va aplica operatorul de uniune, este determinant modul în care se realizează aceasta. Optimizarea cererii trebuie să ia în

considerare ordinea în care facem uniunea tabelelor și numărul tuplurilor care iau parte la "join". De asemenea, este important cu care nod al rețelei începem operațiunea de "join". Dacă datele de la mai multe noduri participă la "join" pentru a executa o interogare, trebuie luat în calcul și costul transmiterii rezultatelor intermediare între noduri.

Optimizarea unei interogări se poate realiza de către SGBD sau de către aplicațiile utilizator. Optimizarea se realizează de către SGBD dacă acesta conține algoritmi interni de optimizare a fiecărei interogări. Deși, o optimizare a SGBD ar fi de dorit, nici un modul de optimizare al SGBD nu este întotdeauna suficient de complet pentru a determina cum este cel mai bine să facă operațiunea de "join" între două tabele stocate la noduri diferite.

În cazul absenței unei optimizări oferite de SGBD, programatorul trebuie să optimizeze fiecare interogare scriind cod pentru operațiunile selecție și uniune realizate între noduri pentru fiecare program de aplicație care necesită așa ceva. În acest caz, programatorul trebuie să țină seama de câțiva factori: dimensiunea tabelor; localizarea tabelor la noduri; disponibilitatea indecșilor; eficiența rețelei; necesitatea posibilității de scriere într-un limbaj procedural a unor secvențe de cod ce nu pot fi scrise folosind doar SQL (un limbaj neprocedural); disponibilitatea datelor (fragmente, replici, snapshot-uri); să folosească, pe cât posibil, cod reutilizabil, pentru ca modificarea acestora să fie facilă.

În plus față de aspectele menționate anterior, se mai pot enumera și alte considerente care pot duce la creșterea complexității problemei și care, cu siguranță, trebuie luate în calcul la optimizarea interogărilor distribuite. Trebuie amintit de aspecte legate de securitatea sistemului și implicațiile autorizării accesului la date: cine ce anume date poate accesa și la ce nod trebuie implementate rutinele de autorizare. Intercalarea rutinelor de autorizare în fluxul de execuție al operațiunilor în sistem duce la creșterea timpului de răspuns a acestuia la interogări.

Un alt aspect este legat de disponibilitatea resurselor. La un moment dat este posibil, și cu

siguranță se întâmplă în decursul timpului ca unul din noduri, să nu fie disponibil pentru interogare, sau chiar rețeaua de comunicație să nu fie utilizabilă. Acest lucru crește timpul de răspuns și implicit costul total, sau, chiar mai rău, poate face imposibilă execuția interogării în cazul în care fragmentul ce conține datele dorite nu este replicat la un alt nod disponibil în acel moment.

Constrângerile de integritate declarate la nivelul bazei de date globale ajută la optimizarea interogărilor adresate sistemului. Se pot astfel elimina din arborele cererii anumite ramuri ce nu se pot întâlni datorită acestor constrângeri de integritate.

Procesul de optimizare a interogărilor distribuite este în strânsă dependență de implementarea și modul de utilizare a rețelei de comunicație. Traficul în rețea poate varia de la o zi la alta, chiar de la o oră la alta, știut fiind faptul că există perioade din zi cu vârfuri de sarcină pentru rețea și perioade cu trafic foarte scăzut. De altfel, pe măsură ce rețeaua suferă modificări, de orice tip, soluția de optimizare trebuie revăzută, și probabil chiar modificată, în cazul în care apar noi posibilități mai rapide de comunicare între noduri.

Am enumerat câteva aspecte cu impact semnificativ asupra timpului de răspuns și implicit al costului total de execuție al unei interogări distribuite. De asemenea, cu cât crește numărul de factori care influențează procesul de execuție al cererilor, cu atât optimizarea acestuia devine o activitate deosebit de complexă.

Concluzii

Optimizarea cererilor în sistemele distribuite este o activitate complexă care depinde de mulți factori. În parte ea este realizată de SGBD, dar există situații în care aplicațiile utilizator trebuie să conțină și algoritmi de optimizare a cererilor.

Dacă fragmentarea bazei de date este realizată corect, o mare parte din cererile utilizatorilor se vor executa la nivel local și astfel execuția cererii va depinde în mare măsură de indecșii pe care i-am creat în deplină concordanță cu predicțiile cererii. Pe de altă parte, politica de securitate a sistemului va restricționa accesul utilizatorilor la o plajă de înre-

gistrări bine precizată, asupra cărora acționează și un set de constrângeri de integritate. Cu toate că putem avea o bază de date bine proiectată astfel încât proporția prelucrărilor locale să fie foarte mare și trierea cererilor de către rutinele de securitate să fie drastică, totuși, în decursul exploatării sistemului distribuit se va dori execuția unor cereri mai complexe care să necesite date de la mai multe noduri. Totodată, în funcție de frecvența de execuție a cererilor în sistemul distribuit, se va justifica o căutare completă în spațiul soluțiilor obținându-se optimul absolut sau se va opta pentru o soluție apropiată de optim găsită cu ajutorul algoritmilor euristici. În aceste cazuri va trebui să ținem seama de aspectele de optimizare menționate în acest articol.

Bibliografie

- [BERN81], P.A.Bernstein, N. Goodman, E. Wong, C.L. Reeve, J.B. Rothnie –*Jr. Query Processing in a System for Distributed Databases (SDD-1)*- ACM Trans. Database Syst., 1981.
- [EPST78], R. Epstein, M. Stonebraker, E.Wong – *Query Processing in a Distributed Relational Database System*- In Proc. ACM SIGMOD Int. Conf. on Management of Data, Austin, 1978.
- [KAMO94], K. Karlaplem, N. Moon Pun - *Query Driven Data Allocation Algorithms for Distributed Database Systems* - RGC CERG, 1994.
- [KOSS00], D. Kossman – *State of the art in distributed query processing* – ACM Computing Surveys, decembrie 2000.
- [LUJA02], G. Lumpkin, H. Jakobsson – *Query Optimization in Oracle9i* – An Oracle White Paper, februarie 2002.
- [LUBO95], I. Lungu, C. Bodea, G. Bădescu, C. Ioniță – *Baze de date, Organizare, Proiectare și implementare* - Ed. ALL, 1995.
- [MULL96], C.S.Mullins – *Distributed query optimization* – Technical Support Magazine, iulie 1996.
- [NAKA94], S.B. Navache, K. Karlapalem, M. Ra – *A Mixed Fragmentation Methodology For Initial Distributed Database Design* – College of Computing, Georgia Institute of Technology, 1994.
- [OSVA95], M. T. Oszu, P. Valduriez – *Principles of distributed database systems* – Prentice-Hall International, 1995.
- [SMBL02], M. Smiljanic, H. Blanken, M. van Keulen, W. Jonker –*Distributed XML Database Systems* – Twente University, Database Group, octombrie 2002.
- [ULLM82], J. D. Ullman - *Principles of database systems* - Rockville, Computer Science Press, 1982.