

SISTEME INFORMATICE PENTRU ASISTAREA DECIZIEI BAZATE
PE SINTEZA DATELOR.
DEPOZITE DE DATE (DATA WAREHOUSE)

Obiective:

- însusirea conceptelor cu privire la sistemele informatice pentru asistarea deciziei bazate pe analiza si sinteza datelor;
- utilizarea tehnologiilor moderne Data Warehousing si On-Line Analytical Processing (OLAP) pentru transformarea datelor în informatii de sinteză;
- însusirea tehnicilor si metodelor de prelucrare multidimensională a datelor.

Concepte cheie: depozite de date (Data Warehouse); prelucrare analitică on-line (OLAP); cubul OLAP; hipercub; bază de date multidimensională. Modul în care datele sunt retransformate în informatii si apoi în cunostinte este de fapt un proces de valorificare a datelor care se realizează prin sintetizarea si analiza lor si în final prin interpretare. Procesul de sintetizare a datelor presupune centralizarea lor, având în vedere diverse criterii si este utilizat în crearea situatiilor de sinteză necesare informării managerilor ca support pentru luarea deciziilor.

Solutiile oferite de informatică pentru procesul de sintetizare a datelor sunt: programe specifice si dedicate; interogări care dau posibilitatea grupării datelor după criterii stabilite si oferă functii pentru domeniile astfel create; functiile de total si subtotal oferite de generatoarele de rapoarte care permit indicarea ierarhiilor criteriilor de grupare.

În ultimul timp, problema centralizării datelor a rămas aceeași, însă volumul de date de explorat este imens, ceea ce duce la faptul că metodele clasice să devină ineficiente. De aceea câștigă tot mai mult teren tehnologii moderne ca Data Warehousing (depozitarea datelor) si OLAP (On-Line Analytical Processing) pe măsură ce suporturile soft devin suport de date pentru sistemele tranzactionale.

Tehnologiile de centralizare transformă datele în informatii de sinteză si analiza lor.

Analiza datelor presupune a găsi relatii între datele sintetizate cum ar fi: asocieri, corelatii structurale, cauzale sau functionale. O formă simplă de analiză a datelor este compararea datelor cu date similare, comparare care se face păstrând toate criteriile identice, doar unul singur având valori diferite.

Compararea se face între seturi de date comparabile, iar tehnologiile de comparatie sunt dotate cu tehnici de observare pentru semnalizarea tiparelor, corelatiilor, asocierilor prin similitudini sau sesizează abateri, exceptii. Informatica a venit în întâmpinarea acestor cerinte cu tehnicile de prezentare grafică care transformă informatia cantitativă de informatie calitativă. Au apărut si tehnici de observare analitică a datelor care au la bază teorii matematice prin care datele reale sunt comparate cu date teoretice produse de un model ipotetic.

Dezvoltarea tehnicilor de observare a dus la aparitia tehnicilor de observare automată bazate pe data-driven. Rezultatul unor astfel de tehnici se regăsesc într-un model cu caracter general. Tehnicile de observare analitică a

datelor se regăesc într-o tehnologie modernă denumită Data Mining (în traducere liberă „Mineritul datelor”).

Rezultatul procesului de observare analitică este obținerea unor tipare, corelații și uneori modele din care se pot deduce tendințe sau se poate prezice cu o anumită probabilitate cum vor arăta datele pe o perioadă ulterioară. Modelul permite interpretarea datelor, ce reprezintă un proces cognitiv cu o apreciere generală a situației, și identifică probleme, oportunități sau potențiale cauze de eșec.

De remarcat este faptul că interpretarea datelor duce la apariția de cunoștințe noi care se vor cumula la cele deja existente. Instrumentele soft clasice pentru asistarea deciziei au avut ca principal scop asigurarea tehnicilor de analiză, optimizare și simulare, precum și reprezentarea grafică a rezultatelor.

Dintre aceste instrumente se amintesc procesoarele de tabele Lotus și Excel orientate pe volume mici de date, cele referitoare la sistemele de gestiune a bazelor de date Access, Visual Foxpro capabile să lucreze cu volume mari de date cu structură uniformă. Principalul dezavantaj al acestor instrumente clasice este că operează numai asupra acelor date care au o structură prestabilită și provin dintr-o sursă unică. Noile sisteme de asistare a deciziei folosesc tehnici speciale de comasare a datelor stocate în structuri neuniforme, pentru a utiliza informații implicite care nu sunt specificate în datele existente. Suporturile software de asistare a deciziei oferă utilizatorilor o serie de facilități cum ar fi: interogarea în limbaj natural, accesul la modele conceptuale, sisteme de gestiune OLAP și servicii de integrare cu alte suporturi soft.

Depozite de date (Data Warehouse)

Necesitatea depozitelor de date este dată de volumul imens de date acumulat în timp de companii. Integrarea acestor date istorice ale companiei într-o structură care să stea la baza luării deciziilor a devenit principala preocupare a noilor tehnologii.

Sistemele de asistare a deciziei care au la bază sinteza și analiza datelor realizează comasarea, sistematizarea, corelarea și gruparea datelor pentru a obține informații care să reliefeze factorii care influențează pozitiv sau negativ performanțele companiei. Ca urmare a obținerii unor astfel de informații se poate adopta o strategie de ameliorare a factorilor cu influență negativă. Obținerea rezultatelor, sub formă de rapoarte care conțin informații utile factorilor de decizie sunt într-o formă accesibilă și sunt rezultatul tehnicilor speciale de explorare a masivelor de date. Aceste tehnici duc la evidențierea unor corelații între date, pot face estimări și prognoze precum și să atragă atenția asupra unor disfuncții.

În sinteză tehnicile de exploatare a masivelor de date pot sugera soluții și pot contribui la luarea deciziilor într-o anumită situație

Datele, mai precis structurile de date care fac obiectul sistemelor informatice de asistare a deciziilor sunt denumite depozite de date (Data Warehouse).

Caracteristicile acestor structuri este faptul că ele pot înmagazina volume mari de date preluate din arhive și/sau din bazele de date ale aplicațiilor informatice specifice activității curente a întreprinderii (sunt volume de ordin

1012 terabytes). Exploatarea acestor volume uriase de date este asigurată de existenta unor motoare speciale care dau posibilitatea ca masivele să poată fi interogate, precum și existenta unor servicii speciale de analiză on-line a datelor (OLAP). Suporturile software sustin performantele prin transformarea datelor, corelarea și completarea lor, precum și prin crearea dictionarului de date, toate acestea asigurând accesul la structurile primare.

Datele sunt extrase din baze de date eterogene create de sistemele informatice deja existente în companie pe diversele platforme hard și soft. Se poate remarca faptul că datele sunt introduse nu la întâmplare, ci sub controlul unor aplicații și al SGBD-ului. Acestea asigură prin serviciile de integritate, stocarea și lucrul în condiții de siguranță maximă. Datele care formează suportul pentru tranzacțiile primare sunt apoi prelucrate pentru a se obține informațiile de sinteză necesare planificării și luării deciziilor și sunt tratate de instrumentele SGBD.

Deoarece exploatarea unui volum enorm de date, pentru a obține diverse rapoarte, este asigurată de integritatea și coerența bazei de date, reuniunea tuturor acestor date duce la exploatarea unui mare număr de tabele, la crearea unor multiple legături virtuale și tabele temporare. Acest volum mare de muncă conduce la principalul inconvenient al depozitelor de date și anume timpul mare necesar exploatarea lor. Un alt inconvenient îl constituie și aglomerarea motorului bazei de date cu task-uri de centralizare care încetinește astfel tranzacțiile curente.

Astfel a apărut necesitatea stocării datelor care sunt dedicate planificării și deciziilor strategice într-un sistem diferit de sistemul operational în așa fel încât funcționarea celor două sisteme să se facă fără inconveniente. În depozitul de date se pot stoca atât arhive de date privind activitatea anterioară, cât și date referitoare la tranzacții ulterioare fără ca utilizatorul să poată interveni.

Datele se pot înmagazina pe domenii sau activități specifice departamentelor unei organizații în așa numitele magazine de date (Data Marts), separarea lor în acest fel ducând la creșterea performanțelor în exploatare. Aceste depozite de date se construiesc de obicei cu tehnologii relationale. Depozitele de date sunt o concentrare de date care organizează, consolidează și centralizează datele din surse eterogene și care vor constitui baza procesărilor analitice atât de necesare proceselor de decizie. Depozitul de date se construiește progresiv adică el permite completări și dezvoltări ulterioare. Pentru a se asigura o calitate sporită a datelor acestea sunt supuse unui proces de curățire și transformare, menționând și maniera de obținere a unor date colectate pe baza celor existente, acest proces ducând la micșorarea timpului cerut pentru obținerea unor rapoarte finale. În depozitele de date se face transformarea codurilor în date explicite, precum și integrarea datelor din nomenclatoare în datele referitoare la tranzacții. Acesta este numit și proces de denormalizare și este caracterizat de faptul că nu modifică integritatea datelor și grăbește procesul de regăsire. Într-un depozit de date redundanța datelor este permisă.

Diferențele dintre depozitul de date și baza de date sunt următoarele
a. Datele continute de un sistem de prelucrare a tranzacțiilor, OLTP (On-Line Transaction Processing) sunt de tip operational, iar datele continute

de un depozit de date sunt specifice asistării deciziilor, sunt date centralizate sau derivate din date operationale, nu se modifică în timp și sunt destinate utilizatorilor finali.

b. În cazul sistemelor tranzactionale, performantele se referă la integritate, confidentialitate, siguranță și timp de răspuns întrucât un număr mare de utilizatori introduc date în sistem, în timp ce în cazul SIAD (deci a depozitelor de date) numărul de utilizatori finali (manageri) este foarte mic.

Astfel și securitatea și siguranța în exploatare nu sunt supuse unor riscuri majore, procedurile de salvare și restaurare fiind mai puțin utilizate decât în cazul sistemelor tranzactionale.

c. Datele procesate în sistemele tranzactionale sunt în seturi relativ mici, introduse recent și compact, astfel încât prelucrarea se face destul de rapid. În procesele decizionale, datele necesare acestora sunt în volum mare, stocate dispersat ceea ce duce la o prelucrare mai lentă.

d. Bazele de date construite pentru sisteme tranzactionale sunt proiectate și realizate pe baza unor cerințe cunoscute și certe, modificările care intervin datorită adaptării sistemului la schimbările intervenite reiau anumite faze ale ciclului de viață. Dar odată implementate ele funcționează perioade lungi de timp fără modificări. În SIAD cerințele sunt cunoscute doar parțial în momentul proiectării și realizării lor, ceea ce obligă depozitul de date să se adapteze din mers cerințelor. De aceea se observă că datele gestionate pentru sisteme tranzactionale sunt privite ca un întreg, pe când cele din depozitele de date sunt organizate pe secțiuni deoarece ele sunt organizate în funcție de subiectul de analiză.

e. Sistemele tranzactionale reflectă de obicei fluxul datelor din activități curente, pe când depozitele de date sunt orientate pe subiecte cum ar fi de exemplu: resurse, produse, clienți, furnizori.

Ciclul de viață al depozitelor de date Depozitul de date (Data Warehouse) este o colecție de date orientate pe subiecte, integrate, corelate în timp și non-volatile care sprijină decizia

Datele care fac obiectul unui depozit sunt integrate în acesta utilizând convenții pentru măsurători, atribute. Structura de care dispune depozitul de date prevede identificarea punctuală a datelor stocate și, mai ales, un acces rapid la ele.

Proiectarea structurii depozitului de date se face prin modelare multidimensională, structura implementându-se ca o bază de date care asigură stocarea unui volum mare de date și un acces rapid la ele, așa numitele baze de date client/server.

Popularea depozitelor de date se face prin preluare din sisteme tranzactionale, dar care vor fi supuse unor procese complexe de transformare care să corespundă structurii depozitului care a fost proiectat. După această etapă, depozitul va putea intra în exploatare pentru a obține analize și rapoarte.

Etapele enumerate anterior (proiectare, populare, exploatare) sunt asistate de un soft specializat de la browsere și generatoare de rapoarte până la instrumente specifice Data Mining.

În exploatarea curentă a depozitului frecvent vor apărea noi cerințe informaționale care vor duce neapărat la extinderea structurii, la popularea cu

extensii cuprinzând date istorice, precum și la integrarea noilor date încorporate în aplicații de analiză. Pe parcursul existenței sale, un depozit de date este incremental și ciclic. Modelarea conceptuală a depozitului de date În etapa de concepție a unui depozit de date se folosesc modele dimensionale care grupează datele din tabelele relationale în scheme de tip stea sau fulg de zăpadă. În aceste scheme pot fi regăsite date cantitative cum ar fi cantități sau valori sau grupate după diverse alte criterii (pe client, pe produs, pe tipuri de servicii etc.). Datele cantitative din bazele de date dimensionale sunt de tip medii, număr de tranzacții, centralizări după anumite caracteristici, totaluri și reprezintă măsuri ale activității. Pe de altă parte, criteriile de agregare vor fi denumite dimensiuni. Măsurile identificate prin dimensiuni vor fi stocate într-un tabel relational care este denumit tabel de fapte, iar codurile utilizate sau asociate criteriilor de agregare sunt date de tabelele de tip nomenclator asociate fiind cu tabelele de fapte și în acest fel schema relatională va fi de tip stea. Dacă se reunesc mai multe scheme de tip stea care utilizează aceleași nomenclatoare formează un model tip constelație. Dacă nomenclatoarele se pot divide în subnomenclatoare atunci există o dependență între acestea. De remarcat că pentru același cod pot exista mai multe nomenclatoare alternative. Dacă se integrează aceste subdimensiuni și dimensiuni alternative, se creează o schemă sub formă de fulg de zăpadă.

Schemele de tip stea, fulg de nea sau constelație sunt modele conceptuale multidimensionale ale depozitelor de date, având ca rol organizarea datelor pe subiecte necesare procesului de decizie. Schema este deschisă (ea se poate modifica pe tot parcursul vieții depozitului de date).

Modul de utilizare a depozitului de date

Depozitele de date conțin structuri unice, integrate și cumulative necesare procesului de decizie. Administratorul depozitului de date are ca principală sarcină stabilirea accesului partajat al categoriilor de manageri prin asigurarea de parole și drepturi de acces. Datele din depozit sunt accesate selectiv de manageri în funcție de necesitățile acestora. În acest fel se creează colecții specializate pe diverse domenii care se numesc magazine de date (Data Marts). Magazinele de date se pot utiliza și ca structuri intermediare pentru colectarea datelor din surse primare și al căror conținut este descărcat periodic în depozitul de date. Depozitele de date pot lua naștere și printr-o stocare exhaustivă a datelor din sistemele tranzacționale în vederea aplicării tehnologiei Data Mining. Utilizarea tehnologiei Data Mining presupune că procesarea datelor se face fără intervenția utilizatorilor, în background, iar rezultatele sunt păstrate pentru a fi consultate ulterior la cerere.

Mediul de depozitare al datelor

Mediul în care se construiește și se exploatează un depozit de date conține următoarele elemente: surse de date tranzacționale, instrumente de proiectare dezvoltare, instrument de extracție și transformare a datelor, sistemul de gestiune al bazei de date, instrumente de acces și analiză a datelor și instrumente de administrare

Toate componentele enumerate sunt integrate pe o platformă Microsoft în mediul de lucru Data Warehousing Framework ca și în cazul SQL Server 7.0.

Acest mediu de lucru oferă asistarea proiectării, implementării și administrării depozitelor de date pe durata vieții (existenței) acestuia. Se poate concluziona că Data Warehousing Framework oferă o arhitectură care se poate integra relativ simplu cu produse ce provin de pe alte platforme, asigură servicii de import-export cu validare și transformarea datelor, asigură metadate integrate pentru proiectarea depozitului și gestionează suportul, task-uri și evenimente. Pentru ca un depozit de date să poată fi procesat este necesară existența unui set specializat de instrumente pentru: descrierea fizică și logică a surselor de date, a depozitelor sau a magaziei de date în care acestea urmează să fie încorporate; validarea, curățirea și transformarea datelor care urmează a fi stocate în depozitul de date; utilizatorii finali, instrumente care permit acestora accesul la datele stocate în depozitul respectiv. Astfel de instrumente sunt specializate pentru medii de dezvoltare a aplicațiilor, produse program specializate pe analiza datelor precum și pentru aplicații personale (individuale). Abordarea multidimensională a datelor stocate 繫 depozite. Definiția și caracterizarea OLAP (On-Line Analytical Processing)

Dacă se analizează tehnologia relatională se observă că cea mai mare parte a problemelor tratate relational sunt în realitate multidimensionale. În modelul relational problemele sunt tratate în tabele care au două dimensiuni: linie și coloană. Problemele reale, care în cea mai mare parte a lor sunt multidimensionale, nu impun limite stocării spațiale a datelor. Astfel, un SGBDR obișnuit nu poate face față cerințelor de agregări de date, sintetizări, consolidări și proiectii multidimensionale. De aceea, a apărut necesitatea extinderii funcționalității unui SGBDR prin adăugarea unor componente speciale care să permită modelare și analiză multidimensională (OLAP) și Data Mining.

Noua tehnologie OLAP permite utilizatorilor navigarea rapidă de la o dimensiune la alta și facilități sporite de obținere a celor mai detaliate informații. Tehnologia OLAP se bazează pe 11 principii formulate de Ted Codd (1992).

Acestea sunt:

- 1) abordarea conceptuală multidimensională a datelor;
- 2) asigurarea unei transparente sporite prin existența unei arhitecturi deschise a sistemului;
- 3) accesibilitatea asigurată utilizatorului prin asistarea implicării acestuia în modalitățile tehnice de furnizare a datelor;
- 4) complexitatea dimensională a analizei oferă performanțe stabile;
- 5) utilizarea arhitecturii client-server, unde server-ul are ca scop omogenizarea datelor;
- 6) posibilitatea de a efectua aceleași operații asupra tuturor dimensiunilor și care poartă numele de prelucrare generică a dimensiunilor;
- 7) gestionarea dinamică a matricilor încrucișate prin facilitatea de a elimina combinațiile dimensionale nule, pentru a nu încărca memoria calculatorului;
- 8) posibilitățile de acces simultan a mai multor utilizatori (multi-user) la aceeași fază (etapă) de analiză;
- 9) operații nerestricțive, ceea ce dă posibilitatea executării fără restricții a calculelor pentru toate combinațiile de dimensiuni și niveluri ierarhice;
- 10) posibilitatea manipulării intuitive a datelor;
- 11) număr nelimitat de niveluri de agregare și de dimensiuni

OLAP este tehnologia de agregare a datelor stocate în depozite într-o manieră de abordare multidimensională cu facilități referitoare la accesul la informații a managerilor în mod interactiv și flexibil. Legătura dintre OLAP și depozitele de date este aceea că OLAP le completează prin transformarea volumului imens de date stocate și gestionat în depozite în informații utile procesului de decizie. Cele 11 reguli ale lui Codd au fost apoi regrupate într-un test cu 5 reguli denumit FASMI (Fast Analysis Shared Multidimensional Information).

OLAP presupune existența unor tehnici care permit de la o navigare și selecție simplă a datelor până la analiza detaliată și complexă. Aplicațiile care se rezolvă pe baza acestei tehnologii au la bază analiza rapidă a informației multidimensională dispersată în locații multiple dar accesibile unui mare număr de utilizatori. Pentru utilizarea acestor facilități, OLAP dispune de eficacitatea bazelor de date multidimensionale și de posibilitatea de a construi alternative pentru diverse probleme de decizie. OLAP presupune că analiza datelor (care pot fi de tip numeric sau statistic) poate fi predefinită de cel care creează aplicația sau chiar de utilizatorul final.

OLAP se caracterizează prin: perspectiva multidimensională a datelor, capacitatea de calcul intensiv și orientare în timp (time intelligence)

Aspectul multidimensional al datelor este dat de posibilitatea de a integra multiplele aspecte care caracterizează activitatea unei întreprinderi și care sunt considerate din perspective multiple ca: timp, bani, produse. Fiecare dimensiune este definită în genere prin mai multe niveluri ca de exemplu: timpul este divizat în an, trimestre, luni, sezoane; produsul în: categorii, clasă. Conceptul de dimensiune este folosit ca înțeles de aspect, dimensiunile fiind independente și cu unități de măsură specifice dimensiunii respective.

Unitățile de măsură pot constitui criterii de agregare a datelor, iar nivelele unei dimensiuni formează ierarhia care la rândul ei poate constitui criteriu de agregare a datelor. Privite din punct de vedere multidimensional, datele sunt reprezentate în hipercuburi de date, prin extinderea cubului tridimensional la cel n-dimensional.

Pe acest tip de cub se pot efectua calcule prin aplicarea unor algoritmi complecși asupra datelor structurate în acesta. Acestea implică posibilitatea de adresare multidimensională directă a cuburilor unitare și optimizarea timpului de răspuns. Caracteristica de orientare în timp (time intelligence) presupune flexibilitatea exploatarea acestei dimensiuni care este necesară pentru comparații și aprecieri de valoare în analizele economice. Această dimensiune este luată de obicei din calendarele tranzacțiilor economice așa cum se află în bazele de date ale sistemului informatic al companiei. Se pot face astfel grupări pe dimensiuni ca: trimestre, luni, ani, sezoane. Se pot utiliza și dimensiuni speciale cum sunt: perioada curentă, perioada precedentă, aceeași perioadă din anul..., care trebuie neapărat luate în considerare la proiectarea hipercubului. Bazele de date multidimensionale folosite de OLAP sunt suprapuse depozitelor de date și stochează straturi de date agregate pe diferite criterii ierarhice. De asemenea, aceste baze de date multidimensionale conțin și date statistice pentru fiecare nivel de agregare.

Modelarea dimensională – cuburi OLAP

Modelarea dimensională presupune conceptualizarea și reprezentarea aspectelor măsurabile ale activității studiate în interdependentă cu contextul în care acesta se desfășoară, aspect identificat prin parametrii activității. Legătura dintre valorile înregistrate ale activității (valori vânzări, cheltuieli comune, costul produselor) și contextul de desfășurare al acestora formează baza numeroaselor rapoarte de sinteză care sunt produse de sistemele tranzacționale. Prin modelare dimensională se oferă un model conceptual comun acestor rapoarte și agregarea lor într-o structură uniformă și flexibilă. Totodată se păstrează și legătura cu sursele inițiale de date, deci posibilitatea de descompunere a datelor centralizate pe niveluri din ce în ce mai mici până se ajunge la setul de tranzacții inițiale (drill-down).

Cubul OLAP se consideră a fi element structural pentru datele din procesul on-line. Acesta este o structură multidimensională, un hiper-cub prin care se modelează complexul de activități pe o perioadă îndelungată de timp. Acest tip de modelare este caracterizat de câteva concepte de bază:

- Cuantificarea activității (aspectul cantitativ) care se face prin utilizarea unităților de măsură clasice ca de exemplu: m, m³, kg, unități monetare. Măsurile cantitative sunt: volum vânzări, volum salarii, cost materiale, cost produs etc.
- Dimensiunile activității sunt de fapt parametrii activității măsurate ca de exemplu: zi, lună, trimestru, client sau grupă de clienți. Dimensiunile sunt de obicei de natură diferită și răspund la întrebări de tipul: Unde? Când? Cu ce? etc.
- Faptele sunt colecții ale cuantificării activității precum și dimensiunile care identifică modul în care acestea s-au desfășurat. Sursa de existență a faptelor este constituită din înregistrările stocate în tabelele de tranzacție ale aplicațiilor operationale care susțin activitatea respectivă. Se pot folosi și dimensiuni scenarii care pot stoca în tabelele de fapte și măsuri imaginare alături de cele reale, pentru ca utilizatorul să poată stoca valori estimate pentru o măsură.

În bazele de date tranzacționale, dimensiunile sunt de fapt câmpuri care conțin caracteristicile unei tranzacții adică datele de identificare ale tranzacțiilor care sunt de obicei chei externe care fac legătura cu nomenclatoarele care le explicitează.

Ca atare, se poate afirma că dimensiunile se materializează în setul de valori posibile care formează domeniul caracteristicii respective, valori care poartă numele de membrii dimensiunii.

O altă caracteristică a dimensiunii este aceea că poate avea mulți adică sunt grupe de valori ale dimensiunii cu o caracteristică comună. Grupele pot fi identificate prin atribute care se află în nomenclatoare și pot lua aceeași valoare pentru mai multe valori ale cheii primare. Multiplii unei dimensiuni nu trebuie să fie neapărat de aceeași natură cu dimensiunea primară, aceasta putând avea mai multe tipuri de multipli în funcție de caracteristicile luate în considerare. Se poate afirma că dimensiunile împreună cu multiplii lor formează structuri arborescente care sunt recunoscute de OLAP ca fiind ierarhii. Ierarhiile pot fi regulate, adică toate ramurile au același număr de ramificații sau neregulate dacă pe anumite ramuri lipsește un nivel de semnificație. La rădăcina arborelui se află o caracteristică cu aceeași valoare pentru toți membrii dimensiunii de bază. Acest tip de caracteristică este una implicită ca, de exemplu, unitatea care are ca activitate cea

analizată sau „all”. Frunzele arborelui formează membrii dimensiunii initiale, iar dimensiunile intermediare pot fi pe mai multe nivele. Dacă arborele este neregulat, pentru a uniformiza ierarhia se poate introduce un membru de tip „alte”.

În acest fel se constată că centralizările pe nivelul respectiv nu vor fi de 100% din valoarea centralizată pe nivelul cel mai de jos. Atributele care definesc ierarhia sunt atribute derivate din atributul care definește dimensiunea acțiunilor măsurate, prin referire la nomenclatoare sau prin clasificări ale valorilor pe care le poate lua atributul respectiv. De exemplu, furnizorii se pot clasifica în furnizori stabili dacă compania face tranzacții cu ei de mai mult de 4 ani, furnizori noi dacă au vechime cuprinsă între 1 și 4 ani și furnizori volatili sau ocazionali dacă în câmpul respectiv din Furnizori nu este completat nimic. Din acest exemplu se observă că asemenea clasificări conduc la obținerea unor atribute derivate prin calcul din caracteristicile aflate în nomenclatoare. În acest fel se vor obține seturi de membri calculați ai dimensiunii. Dimensiunile ierarhizabile se constituie în ierarhii alternative. Nivelele ierarhiilor sunt văzute ca nivel de agregare pentru valorile stocate în tabele de fapte. Membrii dimensiunilor identifică măsura activității stocată în tabelul de fapte. Dacă unui fapt îi sunt asociate mai multe dimensiuni, identificarea unică a acestuia va necesita valori precise pentru fiecare dimensiune. Ca urmare, din tabelele de fapte sunt selectate mai multe înregistrări, adică toate valorile posibile asociate dimensiunilor nespecificate.

Pentru dezvoltarea unui depozit de date, modelarea datelor are un rol important deoarece permite vizualizarea structurii înainte ca ea să fie construită.

Modelul multidimensional reprezentat prin el va fi prezentat desfășurat în secțiuni sau în proiecții tridimensionale.

Secțiunea unui hiper-cub este definită ca o secțiune din cub dată prin coordonatele sale. Proiecția este definită ca o secțiune care centralizează datele de pe toate dimensiunile suprimate.

Vizualizarea on-line se face de fapt tot în secțiuni sau proiecții tridimensionale. Datele din celule sunt prezentate numai în secțiuni sau proiecții transversale bidimensionale. Hiper-cubul ar putea fi imaginat ca un set de tabele pivot grupate pe dimensiunea cerută. Pentru procesul de modelare, hiper-cubul se poate prezenta în formă tabelară în care măsurile sunt evidențiate pe coloane, iar liniile reprezintă combinațiile de dimensiuni. De asemenea, în plan fizic, hiper-cubul poate fi stocat într-un tabel cu coloane multiple în care se stochează măsurile și cu identificatori pe rânduri. Identificatorii de rânduri sunt de fapt chei formate din toate combinațiile posibile de valori ale dimensiunilor.

Utilizarea indecșilor pentru acces rapid nu are prea mare eficiență întrucât cheia este compusă din mai multe caracteristici, iar câmpurile de valoare sunt puține și numerice, astfel că tabelul de indecși este aproape de aceeași dimensiune cu tabelul inițial. De aceea, se utilizează tabelul bitmap pentru un acces direct rapid. Datele modelate ca hiper-cuburi formează baze de date multidimensionale.

Baze de date multidimensionale

Baza de date multidimensională este formată din două structuri:

structura datelor în care se stochează măsurile activităților preluate din tabela

de fapte a depozitului de date. Datele vor fi prezentate utilizatorului în celulele tabelor pivot; structura metadatelor care este formată din totalitatea dimensiunilor și membrilor acestora precum și din structurile ierarhice ale dimensiunilor. Utilizatorul poate vizualiza această structură ca nume de coloane și linii care reprezintă informațiile de pe axele cuburilor. Numerotarea nivelurilor începe de la rădăcină (nivel 0) către frunze (unde va apare nivelul maxim). Ierarhiile posedă propriile lor seturi de niveluri, chiar dacă unele ramuri sunt comune. De exemplu: ierarhia Calendar este formată din nivelele (0-5): Timp, An, Semestru, Trimestru, Lună, Dată calendaristică, ierarhia Anotimp este formată din nivelele (0-4): Timp, An, Sezon, Lună, Dată calendaristică, iar ierarhia Anotimp este formată din nivelele (0-3): Timp, Săptămână, Zi, Dată calendaristică. Pe fiecare nivel se stochează membrii dimensiunilor respective. Rădăcina care se observă că este comună (Timp) este nivelul de agregare maxim având ca unic membru implicit „all”. Orice nod în arbore este un membru al unei subdimensiuni. Nodurile subordonate unui nod formează un set, iar orice membru al unui set are un număr de ordine începând cu 0. De asemenea, orice membru poate avea proprietăți ca de exemplu unele zile sunt sărbători legale, unii ani sunt bisecti. Exemplul prezentat presupune o structură strict arborescentă întrucât fiecare membru al unei dimensiuni are submembri distincti, chiar dacă aceștia au aceleași valori. De exemplu, fiecare an are setul lui de luni, fiecare săptămână are setul ei de zile. Ca mod de identificare, membrii vor fi calificați cu numele membrului de pe nivelul precedent căruia acesta i se subordonează: 2000-feb, 2001-feb. Tipul acesta de dimensiuni care au membri ce se repetă se pot crea și ulterior prin combinarea a două nivele din ierarhie sau din ierarhii diferite pentru a crea un nivel nou, virtual.

Pentru a se putea naviga pe o structură arborescentă, sistemele de gestiune pun la dispoziție operatori ierarhici. De exemplu, pentru exploatarea datelor, sistemele de gestiune oferă operatori pe hipercuburi. Fizic, datele sunt stocate într-un fișier cu acces direct pe baza adresei fizice absolute sau relative a înregistrării obținute prin exploatarea tabelor bitmap obținute în urma creării structurii de date. Aceste tabele sunt puntea de legătură dintre structura de date și structura de metadata. Iată cum se face această legătură: se știe că pentru fiecare membru al fiecărei dimensiuni există o coloană (1 bit) în tabele bitmap pentru fiecare înregistrare există un rând în același tabel în care se stochează 1 în dreptul biturilor asociați membrilor dimensiunii existente în înregistrare. Datorită acestui procedeu, câmpul respectiv nu trebuie stocat în înregistrare, iar structura datelor este redusă la un minim necesar. Din tabelul de măsuri se vor putea selecta acele înregistrări care au un bit 1 în poziția corespunzătoare biturilor 1 din mască. Un inconvenient al tabelor bitmap este acela că ele sunt greu de obținut, iar apariția unor noi membri sunt greu de inserat în poziția corespunzătoare. Procesul de refacere a unui tabel bitmap este mare consumator de timp având în vedere că tabelul de fapte din depozit (care se va transforma în baza multidimensională) poate avea un număr imens de înregistrări. Mască de interogare se obține prin exploatarea structurii ierarhice a metadatelor de unde se pot extrage seturi de membri pentru dimensiunile

desemnate prin specificatorii de axe. Adresarea tabelului de măsuri se face în mod direct pe baza unui set de adrese de înregistrări care se suprapun cu tiparul măştii. Din tabel se preiau în această manieră valorile care se centralizează pentru celula cubului cu dimensiunile sale.

Se poate afirma că structura metadatelor este de tip ierarhic, fiecare dimensiune fiind stocată într-o structură arborescentă cu o singură rădăcină (all) și cu o multitudine de ramuri care pot conține frunze comune (ierarhii alternative). Orice nivel al unei ierarhii poartă un nume și conține un set de membri. De altfel și ierarhiile alternative poartă un nume pentru a putea fi distinse. Structura în care sunt stocate datele este o structură cu acces direct prin tabele bitmap exploatate prin măști.

Operatii OLAP asupra hipercubului

Un hipercub este proiectat astfel încât el să aibă în vedere nivelul de detaliu necesar în procesul de analiză. Nivelul de detaliu (granularitatea) reprezintă numărul de membri ai unei dimensiuni. Datele pot fi vizualizate printr-o selecție în hipercub pe baza unui criteriu ierarhic care ar putea fi de exemplu structura organizatională pe care o conduce un anumit manager. Dacă de la pornire, granularitatea este prea mare, datele vor fi mult prea centralizate și nu se va putea face decât o analiză grosieră. Ajustarea nivelului de granularitate este realizată de OLAP prin exploatarea ierarhiilor dimensiunilor prin comasări și descompuneri ale măsurilor prin proceduri care poartă numele de drill-up și drill-down. Prin intermediul acestor proceduri se face o deplasare a proiecției cubului în sus sau jos pe nivelele ierarhice ale fiecărei dimensiuni (zoom in; zoom out), executând de fiecare dată centralizări ale măsurilor stocate la cea mai mică granularitate după criteriile ierarhice stabilite în prealabil.

Este stabilit un nivel de granularitate inițial sub care nu se poate coborî. Din acest motiv este important ca dimensiunile de bază să fie cât mai rafinate sau să se creeze Data Marts, unde hipercuburile sunt proiectate la nivelul de detaliu stabilit de managementul operational. Pentru managementul superior se va construi un depozit cu hipercuburi centralizatoare cu granularitate mare. Prin drill-down se obțin detalii, iar prin drill-up se obțin date sintetice.

Un alt grup de operații oferit de OLAP este sectionarea (slicing) și defalcarea (dicing). Prin sectionare, se creează posibilitatea selectării prin vizualizare doar pentru un membru al unei dimensiuni, adică un plan din cubul tridimensional. Secțiunea astfel obținută va apărea ca un tabel pilot cu valorile dimensiunilor pe laturi și cu specificarea valorii alese pentru dimensiunea suprimată. Defalcarea (dicing) este operația de proiectare a unei dimensiuni pe o altă. De obicei o dimensiune din primul plan este combinată cu o altă dimensiune din adâncime. Acest proces se mai numește imbricarea dimensiunilor.

Dimensiunile unui cub pot fi private sau pot fi utilizate în comun și de alte cuburi (ele provin din depozitele cu schema de tip constelație). Proiectarea structurilor depozitelor de date și a cuburilor OLAP este un proces ce se desfășoară continuu pe tot parcursul existenței (vietii) aplicației, dimensiunile cuburilor fiind în strânsă dependentă cu detaliile activității structurate.

Aplicatiile construite cu tehnologia OLAP își găsesc locul în multiplele domenii ale activității întreprinderilor, de la finanțe, bănci, marketing până la producție și vânzări. De exemplu, activitatea de producție poate fi susținută de aplicații OLAP cum sunt: planificarea operațiilor, controlul calității produselor, analiza rebuturilor, analiza optimizării raportului dintre cost-beneficii. OLAP, utilizând tehnici inteligente de optimizare, beneficiază de avantajul timpului de răspuns mic.

Crearea aplicațiilor OLAP în Microsoft SQL Server

Pentru realizarea unei aplicații OLAP sunt necesare următoarele etape:

1. Crearea bazei de date relationale (tranzactionale) care va conține datele curente ale organizației rezultate din tranzacții.
2. Crearea bazei de date multidimensionale, a cuburilor și tabelelor de fapte care preiau datele din baza de date relatională. Datele sunt extrase, transformate și încărcate în tabelele de fapte din tabelele relationale.
3. Crearea interfeței aplicației într-un mediu de programare vizual – Visual Basic.

Crearea bazei de date tranzactionale în Microsoft SQL Server

Datele stocate în cadrul organizației sunt importate într-o nouă bază de date tranzactională ce stă la baza construirii cuburilor de date. Datele sunt organizate în tabele care corespund dimensiunilor, ierarhiilor și tabelelor de fapte ale cuburilor multidimensionale.

Ca exemplu, se va crea o aplicație destinată analizei rezultatelor financiare ale unei bănci comerciale. Se vor analiza volumul depozitelor și volumul creditelor în funcție de următoarele dimensiuni: agentie, durată, garanție, monedă, sector de activitate, timp, tip depozit, tip client (pers fizică sau juridică). Tabelul de fapte construit va conține două măsuri: volumul depozitelor și volumul creditelor.

Crearea tabelelor în Microsoft SQL Server se realizează prin utilizarea de scripturi, așa cum se prezintă în exemplul de mai jos:

- create table agentie (Agentie varchar(20), Zona varchar(10), Tara varchar(10));
- create table voldepozite (Agentie varchar(20), Durata varchar(20), Moneda varchar(20), Tipjur varchar(30), Timp varchar(20), Tipdepozit varchar(30), voldep numeric);
- create table volcredite (Agentie varchar(20), Durata varchar(20), Moneda varchar(20), Tipjur varchar(30), Timp varchar(20), Garantiecredit varchar(20), Sectoractivitate varchar(20), volcredite numeric).

Crearea bazei de date multidimensionale în SQL Server Produsul Microsoft SQL Server oferă suportul și instrumentele necesare dezvoltării sistemelor OLAP prin setul de aplicații SQL OLAP Services, iar gestiunea bazei de date multidimensionale este realizată de serverul OLAP.

Se creează o nouă bază de date multidimensională care va conține cuburile de date prin intermediul meniului New Database. Se creează cuburile cu ajutorul asistentului Cube Wizard.

Cuburile OLAP utilizează datele stocate în tabelele bazei de date tranzactionale. Din acest motiv trebuie configurată conexiunea dintre cubul OLAP și baza de date tranzactională din care vor fi preluate datele. Conexiunea cu serverul de baze de date Microsoft SQL Server se realizează cu ajutorul

Microsoft OLE DB Provider for SQL Server. După stabilirea serverului tranzacțional se selectează și baza de date tranzacțională din care se importă datele.

Aplicatia OLAP conține două cuburi pe care se vor analiza cele două tipuri de operațiuni bancare: operațiunile pasive (constituirea de depozite) și operațiunile active (acordarea de credite). Pentru fiecare cub se definește o schemă care conține în centru tabelul de fapte legat de dimensiunile corespunzătoare fiecărei activități analizate după cum urmează :

1. Cubul Depozite – urmărește analiza depozitelor și a dobânzilor pasive rezultate din activitatea curentă a băncii.

- Dimensiunile identificate în cadrul acestei scheme sunt: Agentie, Durata, Moneda, Timp, Tip juridic, TipDepozit.
- Tabelul de fapte al modelului este VolDepozite având ca măsură volumul depozitelor constituite (voldep).

2. Cubul Credite – urmărește analiza creditelor și a dobânzilor active rezultate din activitatea curentă a băncii. În cadrul acestei scheme se identifică dimensiuni comune cu schema operațiunilor pasive. Acestea sunt:

- Dimensiunile identificate în cadrul acestei scheme sunt: Agentie, Durata, Moneda, Timp, Tip juridic, Garantie credit, Sector activitate.
- Tabelul de fapte al modelului este VolCredite ce conține măsura VolCred (volumul creditelor).

Crearea cuburilor se poate face independent (în fereastra Cube Editor) sau în mod asistat (prin opțiunea Cube Wizard).

În ambele cazuri trebuie parcurși următorii pași:

1. Alegerea tabelului din baza de date relațională din care se importă datele. După configurarea conexiunii dintre serverul OLAP și baza de date tranzacțională se pot vizualiza tabelele acestuia. O parte din aceste tabele conțin date necesare dimensiunilor din cuburi, iar altele furnizează date în tabelele de fapte.
2. Stabilirea tabelului de fapte. În cazul cubului Depozite, tabelul de fapte Vol Depozite importă datele din tabelul voldepozite din baza de date relațională, iar în cazul cubului Credite tabelul de fapte VolCredite importă datele din tabelul volcredite din baza de date relațională.
3. Crearea și configurarea dimensiunilor. Se aleg tabelele din baza de date relațională din care se vor importa datele în tabelele dimensiuni.
4. Pentru fiecare dimensiune se stabilesc nivelurile ierarhice. De exemplu, pentru dimensiunea Timp, prezentă în ambele cuburi nivelele ierarhice sunt: Luna, Semestru, An.
5. Alegerea dimensiunilor care intră în configurația cubului. După crearea tuturor dimensiunilor, se aleg numai cele care participă la configurația cubului.
6. Stabilirea opțiunilor de stocare și procesare a datelor din cubul OLAP.