

MAGAZIILE DE DATE – SUPTOR INFORMATIC PENTRU LUAREA DECIZIILOR

Dr.hab., conf. Univ. V. Cotelea
Academia de Studii Economice din Moldova

1. DESCRIERE GENERALĂ

Odată cu apariția calculatoarelor, organizațiile utilizează datele recuperate din sistemele informatice pentru a-și satisface necesitățile informaționale. Unele software-uri oferă accesul direct la datele conținute în sistemele operaționale. Altele extrag datele din bazele de date pentru a le combina în diverse forme nestructurate, ținând seama de necesitățile în informații ale utilizatorilor. Ambele metode au evoluat pe parcursul timpului și acum organizațiile gestionează datele necurățate și inconsistente, în baza cărora, în majoritatea cazurilor, se iau decizii importante.

Este recunoscut faptul că, pentru o creștere a eficienței întreprinderii, este nevoie de cea mai bună utilizare a resurselor informatice care există în interiorul și exteriorul ei. Cu toate că sistemele informatice au început să fie implementate cu mulți ani în urmă, încă nu se poate spune că există o utilizare eficientă a acestora.

Motivul principal îl constituie modul în care au evoluat calculatoarele, tehnologiile informaționale și sistemele. Majoritatea organizațiilor tind să obțină informații calitative, dar atingerea acestui obiectiv depinde de arhitectura sistemelor, atât a hardware-ului, cât și software-ului.

Sistemele care susțin procesele de gestiune și luare a deciziilor sunt esențial diferite de sistemele convenționale, de prelucrare a tranzacțiilor sau care se mai numesc sisteme operaționale, adică, aplicațiile tipice, într-o întreprindere, pot fi clasificate în două mulțimi mari:

- aplicațiile de business, aplicațiile ce susțin businessul de fiecare zi, care garantează funcționarea întreprinderii;
- aplicațiile despre business, aplicațiile ce analizează businessul, care ajută la interpretarea a ceea ce se va întâmpla și la luarea deciziilor asupra strategiilor de viitor.

O arhitectură de date adecvată, care ar suporta ambele tipuri de aplicații, se bazează pe două medii de baze de date: bazele de date operaționale (pentru a susține aplicațiile de business) și bazele de date pentru asistarea

deciziilor (pentru a susține aplicațiile despre business).

Înmagazinările de date (*data warehouses*), actualmente sunt, centrul de atenție al întreprinderilor și instituțiilor mari, deoarece ele asigură un mediu în care sunt mai bine utilizate datele gestionate de diverse aplicații operaționale.

O înmagazinare de date este o colecție de date, în care datele organizației sunt integrate și care se utilizează în calitate de suport pentru procesul de luare a deciziilor manageriale. Deși diverse organizații și persoane fizice au reușit să înțeleagă tehnicile unei înmagazinări de date, experiența a demonstrat că există încă multe dificultăți în această direcție.

Evident că reunirea elementelor de date corespunzătoare, din diverse izvoare de aplicare într-un mediu integral centralizat, simplifică problema de acces la date și, în consecință, accelerează procesul de analiză și consultare.

Aplicațiile pentru susținerea deciziilor bazate pe o înmagazinare de date pot face mai practică și mai ușoară exploatarea datelor și o mai înaltă eficiență a afacerilor, care nu se obține când se utilizează numai datele provenite din aplicațiile operaționale, în care datele se obțin, realizând procese independente și, de multe ori, mai complexe.

Magaziile de date se creează pe datele extrase dintr-una sau mai multe baze de date orientate spre aplicațiile operaționale. Datele extrase sunt transformate (cu scopul de a fi eliminate inconsistențele) și generalizate, dacă este nevoie, și apoi încărcate în înmagazinarea de date. Procesele de transformare, creare a variantelor în timp, generalizare, agregare și combinare a datelor extrase, asigură crearea unui mediu pentru accesul la datele organizației. Această tehnică nouă dă persoanelor fizice posibilitatea de acces la toate nivelele întreprinderii, de a lua decizii cu o mai mare siguranță și responsabilitate.

Noile tehnologii informaționale, în cadrul unui mediu de înmagazinare, oferă organizației posibilitatea de folosire optimală a datelor în calitate de ingredient-cheie pentru un proces eficient de luare a deciziilor. Organizațiile pot scoate

avantaje din resursele lor de informații în realizarea operațiilor de afaceri, dar, pentru aceasta, trebuie să fie considerate strategiile tehnologice speciale de implementare a unei arhitecturi complete a unei magazii de date [1].

2. CONCEPTUL „MAGAZIE DE DATE”

O magazie de date este, în general, o bază de date dedicată. Ea înglobează majoritatea datelor din sistemul operațional al întreprinderii și este separată de operațiile acestuia. Înmagazinarea de date facilitează utilizarea aplicațiilor, organizarea și păstrarea datelor necesare pentru prelucrarea analitică, prelucrarea datelor pe o perspectivă mai largă de timp.

O magazie de date se poate caracteriza prin modul în care datele de afaceri păstrate în înmagazinarea de date diferă de datele operaționale utilizate în aplicațiile de producție (figura 1).

Baza de date operaționale	Magazia de date
Date operaționale	Date de afaceri pentru informare
Orientate pe aplicație	Orientate pe subiect
Actuale	Actuale + istorice
Detaliate	Detaliate + agregate
Schimbare continuă	Stabile

Figura 1. Deosebirea tipurilor de date.

Evident că datele, într-o magazie de date, vin, în majoritatea cazurilor, din mediul operațional. Înmagazinarea de date este întotdeauna o stocare de date transformate și separate fizic de datele mediului operațional.

William H. Inmon, pionier în elaborarea magaziiilor de date, emite următoarea definiție a magaziei de date [2]:

Definiția 1. *O magazie de date* este o colecție de date orientată spre subiect, integrată, variabilă în timp și nevolatilă, care are drept obiectiv susținerea procesului de luare a deciziilor.

Această definiție are două supoziții implicite: un data warehouse este separat fizic de sistemul operațional și susține date agregate și tranzacții de date (atomice), care sunt separate de baza de date utilizată în regim OLTP.

Astfel, o înmagazinare de date reprezintă o bază de date ce conține date extrase din mediul de producție al întreprinderii, care sunt selectate și epurate, optimizate pentru procesul de consultare, și nu pentru procesul de tranzacții. Magaziile de date

presupun consolidarea datelor din mai multe surse, fie cele depozitate în bazele de date relaționale, fie datele provenite din tabelele electronice, documente textuale etc.

În acord cu alt specialist în domeniul înmagazinării de date, Richard Hacathom, obiectivul unei magazii de date constă în “formarea unei imagini unice a realității de business”:

Definiția 2....sistemele data warehouse cuprind o mulțime de programe ce extrag datele din mediul operațional al întreprinderii, o bază de date care le păstrează și sisteme ce furnizează aceste date utilizatorilor săi.

Există numeroase definiții ale magaziiilor de date, primele concentrându-se asupra caracteristicilor datelor păstrate în magazia de date. Definițiile alternative extind domeniul de valabilitate al definiției înmagazinării datelor, pentru a include prelucrarea asociată accesării datelor, de la resursele inițiale până la livrarea acestora către organele de decizie.

Oricare ar fi definiția, scopul suprem al înmagazinării datelor rezidă în integrarea datelor generale din întreaga întreprindere într-o singură magazie, de la care utilizatorii pot lansa interogări, elabora rapoarte și efectua analize. Magazia de date reprezintă un mediu de susținere a deciziilor, care preia datele stocate în diverse surse operaționale, le organizează și le face disponibile organelor de decizie din cadrul întreprinderii. Pe scurt, o magazie de date constituie o tehnologie de administrare și analiză a datelor.

Precum s-a afirmat, principalele caracteristici ale unei magazii de date sunt:

- orientarea spre subiect (temă);
- datele sunt integrate;
- datele sunt istorice, adică variabile în timp;
- datele sunt nevolatile.

2.1. Orientarea spre temă

Prima caracteristică a magaziei de date constă în faptul că datele ei se clasifică, ținând cont de aspectele ce prezintă interes pentru întreprindere, adică, datele sunt organizate pe teme.

Mediul operațional este proiectat în jurul aplicațiilor și funcțiilor, ca prestări, economii și depozite într-o instituție bancară. De exemplu, o aplicație de intrare a ordinelor poate accesa date despre clienți, articole și conturi. Baza de date îmbină aceste elemente într-o structură ce satisface necesitățile aplicației.

Mediul unui data warehouse este organizat în jurul unor subiecte precum client, vânzător, articol și activitate. De exemplu, pentru un fabricant,

acestea pot fi clienții, articolele, furnizorii sau vânzătorii. Pentru o universitate, pot fi studenții, cursurile sau profesorii. Pentru un spital, pot fi pacienții, personalul medical, medicamentele etc. Alinierea în jurul temelor afectează proiectarea și implementarea datelor găsite în magazia de date.

Aplicațiile sunt asociate cu proiectarea bazei de date și a proceselor. Înmagazinarea datelor se focalizează pe modelarea datelor și proiectarea bazei de date, și proiectarea proceselor (în forma sa clasică) nu este separată de acest mediu.

Diferența dintre orientarea pe procese și funcții (ale aplicațiilor) de orientarea pe teme constă în conținutul datelor la nivelul detaliat. Într-o magazie de date, se exclud datele ce nu vor fi utilizate de sistemele de susținere a deciziilor, în timp ce datele orientate pe aplicații sunt datele îndreptate spre satisfacerea imediată a cerințelor funcționale și a proceselor, care pot fi utilizate (sau nu) și pentru susținerea deciziilor.

Altă deosebire importantă este legătura dintre date. Datele operaționale mențin o legătură continuă între două sau mai multe relații bazate pe o regulă comercială aflată în vigoare. Înmagazinările de date măsoară o perioadă de timp în care sunt multe legături. Multe din regulile comerciale (și legăturile corespunzătoare de date) se reprezintă într-o magazie de date între două sau mai multe relații.

Se poate menționa că interesul organizării pe teme constă în faptul că devine posibilă realizarea analizelor pe subiecte, trecând prin structurile funcționale și organizaționale ale întreprinderii. Această orientare permite, de asemenea, petrecerea analizelor prin iterație, subiect după subiect.

Integrarea într-o structură unică este indisolubilă de evitarea duplicatele datelor ce se referă la mai multe subiecte. Uneori, în practică, există și **piețe de date** (*datamarts*), adică, magazia de date este fragmentată în mai multe baze, care suportă orientarea spre temă.

2.2. Date integrate

Cea mai importantă caracteristică a mediului de înmagazinare a datelor constă în faptul că datele dintr-o magazie de date sunt integrate întotdeauna.

Integrarea datelor se prezintă prin multe aspecte: convenții de nume consistente, unități de măsură uniformă a variabilelor, codificări de structuri consistente, atribute fizice ale datelor consistente, izvoare etc.

Pe parcursul anilor, proiectanții diferitelor sisteme luau decizii specifice asupra modului în care trebuie să fie construită o aplicație. Stilurile și proiectele personalizate pot fi reflectate în mai

multe moduri. Ele se deosebesc prin codificare, prin structurile cheilor, prin caracteristicile lor fizice, prin convențiile de denumiri etc. Posibilitatea unui grup de proiectanți de aplicații de a crea aplicații inconsistente este fabuloasă. În continuare, sunt prezentate cele mai importante deosebiri ce pot fi în formele în care se proiectează aplicațiile.

2.2.1. Codificarea

Proiectanții de aplicații pot codifica atributul *Sex* în diverse moduri. Unii proiectanți reprezintă atributul *Sex* prin valorile “M” și “F”, alții prin “I” și “O”, alții prin “X” și “Y”, și inclusiv prin “masculin” și “feminin”.

Nu importă cum va sosi atributul *Sex* în magazia de date. Probabil, valorile “M” și “F” sunt cele mai potrivite. E important faptul că atunci când atributul *Sex* vine din diverse izvoare, trebuie să sosească în magazie în stare uniform integrată. Prin urmare, când atributul *Sex* este încărcat într-o magazie de date dintr-o aplicație și în magazie se păstrează în formatul “M” și “F”, datele trebuie convertite în formatul respectiv.

2.2.2. Măsurarea atributelor

Proiectanții folosesc unități de măsură diferite pentru a prezenta, de exemplu, mărimea țevelor. Unii proiectanți păstrează datele despre țevi în centimetri, alții în inche-uri, alții în milioane de picioare, cuburi pe secundă, iar alții în yarzi.

La definirea unităților de măsură a atributelor, transformarea respectivă traduce diverse unități de măsură utilizate în diferite baze de date într-o măsură standard comună. Astfel, oricare ar fi izvorul când datele despre țevi sunt trecute într-o magazie de date trebuie să fie măsurate în același mod, cu aceeași unitate de măsură.

2.2.3. Convenții de denumire

Același element, în diverse aplicații, frecvent, este referit prin diferite nume. Procesul de transformare asigură faptul că se va utiliza numele preferabil utilizatorului. De exemplu, același concept *balanță*, într-o aplicație, poate fi numit *balanță activă*, în alta *balanță curentă*, în a treia *flux în casă* etc.

2.2.4. Izvoare multiple

Același element poate fi obținut din mai multe surse. În acest caz, procesul de transformare trebuie să asigure că este utilizat izvorul potrivit, documentat și motivat pentru a fi depozitat în magazie.

După cum urmează din exemplele aduse, integrarea afectează aproape toate aspectele proiectului - caracteristicile fizice ale datelor, alternativa de a avea mai mult de o sursă de date, problema de standardizare a denumirilor inconsistente, formate de date inconsistente etc.

Oricare ar fi forma de proiectare, rezultatul trebuie să fie același – datele trebuie să fie păstrate în magazia de date, într-un model global acceptabil și singular, chiar dacă sistemele operaționale subiacente păstrează datele în mod diferit. Or, când analistul sistemului de asistare a deciziilor lucrează cu magazia de date, el trebuie să se concentreze asupra utilizării datelor ce se găsesc depozitate, dar nu să fie preocupat de confiabilitatea și consistența datelor.

2.3. Date istorice sau variabile în timp

Datele istorice sunt necesare pentru a urmări în timp evoluția diferitelor valori ale indicatorilor supuși analizei. Astfel, un timp referențial trebuie să fie asociat cu datele pentru a permite identificarea pe durata valorilor precise.

Datele asociate unui anumit moment sunt solicitate dintr-o magazie de date. Această principială caracteristică a datelor este foarte diferită de cea a datelor din mediul operațional. În mediul operațional, sunt cerute datele asociate momentului de accedere. Cu alte cuvinte, în mediul operațional, când se accede la unitatea de date, se speră că valorile cerute reflectă starea actuală a domeniului de interes.

Deoarece, într-o magazie de date, datele solicitate pot fi asociate oricărui moment de timp (adică, nu “chiar acum”), datele din magazie se numesc *variabile în timp*.

Datele istorice sunt puțin utilizate în procesarea operațională. Datele din magazie, în schimb, trebuie să includă date istorice pentru a fi utilizate la identificarea și evaluarea tendințelor.

Datele istorice sau variabile în timp sunt caracterizate în diverse moduri:

- În primul rând, datele sunt păstrate o perioadă îndelungată de timp – de la cinci la zece ani. Perioada de timp reprezentată în mediul operațional este mult mai scurtă – de la valori actuale la șazeci - nouăzeci de zile.

Aplicațiile care au performanțe bune și sunt disponibile pentru procesarea tranzacțiilor trebuie să aducă o cantitate minimă de date și să posede un oarecare grad de flexibilitate. Datorită proiectării rigide a aplicației, aplicațiile operaționale au un termen de viață scurt.

- În al doilea rând, într-o magazie de date, datele sunt variabile în timp datorită structurii cheii. Orice cheie într-o magazie de date conține, explicit sau implicit, un element de timp cum ar fi ziua, săptămâna, luna etc.

Elementul care reprezintă timpul este, aproape întotdeauna, parte componentă a cheii concatenate. Ocazional, timpul poate exista implicit, ca în cazul când fișierul complet este duplicat la sfârșitul lunii sau trimestrial.

- În al treilea rând, datele unei magazii de date, odată înregistrate corect, nu mai pot fi actualizate. Pentru toate scopurile practice, datele unei magazii de date reprezintă o serie largă de viziuni instantanee (snapshots).

Bineînțeles, dacă snapshoturile datelor sunt luate incorect, atunci pot exista schimbări. Dacă snapshoturile sunt făcute adecvat, ele nu sunt modificate niciodată. În unele cazuri, este neetică și chiar ilegală modificarea snapshoturilor din magazie. În schimb, datele operaționale care trebuie să reprezinte realitatea la momentul accesării, sunt actualizate în acord cu schimbarea domeniului de interes.

2.4. Date nevolatile

Datele, într-o magazie de date, sunt nevolatile, deoarece datele nu sunt reactualizate în timp real, ci sunt reîmprospătate de către sistemele operaționale în intervale regulate de timp. Datele noi sunt adăugate întotdeauna, mai degrabă, ca un supliment al bazei de date, decât ca o înlocuire. Baza de date absoarbe continuu aceste date noi, integrându-le pe rând cu datele anterioare. Cu alte cuvinte, aceeași cerere lansată la diferite momente de timp, asupra acelorași date întotdeauna trebuie să returneze aceleași rezultate.

Astfel, datele într-o magazie sunt utile numai când sunt stabile, adică datele nu sunt șterse. Perspectiva mai mare, esențială pentru o analiză și luare a deciziilor, necesită o bază de date stabilă.

Datele din sistemele operaționale se schimbă din moment în moment, adică actualizarea (inserarea, ștergerea și modificarea) în mediul operațional se face regulamentar, înregistrare cu înregistrare. În schimb, manipularea datelor într-o magazie de date este mult mai simplă. Există două tipuri de operații: încărcarea inițială a datelor și accesarea acestora. Nu există actualizări de date (în sensul general de actualizare) în magazie.

Există unele consecințe ale acestei deosebiri importante între prelucrarea operațională și prelucrarea în înmagazinările de date. La nivelul de proiectare în magazii, nu trebuie ținut cont de

anomaliile de actualizare, deoarece nu are loc actualizarea datelor. Aceasta înseamnă că, la nivelul fizic de proiectare, se poate optimiza accesul la date, în particular, cu utilizarea normalizării și denormalizării fizice.

O altă consecință a simplității operației înmagazinării de date constituie tehnologia subiacentă, utilizată pentru a trimite datele în magazia de date. Suportarea actualizării înregistrare cu înregistrare în mod on line (cum este frecvent în cazul prelucrării operaționale) cere ca tehnologia să poseze un fundament mai complex, dar sub o față simplă.

Sursa aproape a tuturor datelor unei magazii de date este mediul operațional. La o simplă viziune, se poate crede că există o redundanță masivă de date între datele mediilor. Bineînțeles, prima impresie a multor persoane se concentrează spre redundanța mare de date între mediul operațional și mediul unei magazii de date. Această concluzie este superficială și demonstrează o carență în înțelegerea a ceea ce apare într-o bază de date. De fapt, există o redundanță minimă de date între ambele medii, dacă se ține cont de următoarele:

- Datele sunt filtrate când trec din mediul operațional în magazie. Există multe date ce nu vor ieși din mediul operațional. Numai datele ce realmente sunt necesare vor intra în magazia de date.
- Termenul de păstrare a datelor este foarte diferit de la un mediu la altul. Datele din mediul operațional sunt mai recente decât cele din magazia de date. Din perspectivele de timp unic, există o suprapoziție între mediile operațional și magazia de date.

- Magazia de date conține un rezumat de date, care nu se găsește în mediul operațional.

- Datele suportă o transformare serioasă când trec în magazia de date. Cea mai mare parte din date se schimbă semnificativ când sunt selectate și mutate în magazia de date. Altfel spus, majoritatea datelor se modifică fizic când se mută în magazie. Din punct de vedere al integrării nu sunt aceleași date care rezidă în mediul operațional.

Din perspectiva acestor factori, redundanța datelor între cele două medii este o stare rară.

3. STRUCTURA UNEI MAGAZII DE DATE

Magaziile de date dispun de o structură distinctă. Există diverse niveluri de sinteză și detalieri ce delimitează magazia de date. Structura unei magazii de date este prezentată în figura 2.

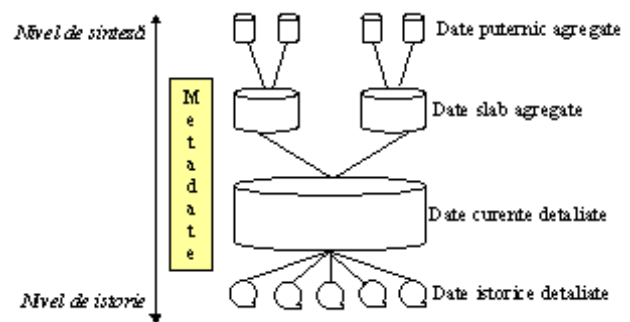


Figura 2. Structura datelor într-o magazie de date

O magazie de date conține cinci clase de date, organizate conform unei axe de istorie și unei axe de sinteză:

- datele curente detaliate;
- datele istorice (mai vechi) detaliate;
- datele slab agregate;
- datele puternic agregate;
- metadatele.

O problemă-cheie ce trebuie soluționată în proiectarea unei magazii de date constă în granularitatea datelor. Datele primitive, precum cele păstrate în sistemele de prelucrare a tranzacțiilor, de obicei, sunt cele mai detaliate.

Agregarea sau sumarea este un proces de intrare a datelor primitive sau detalii de date pentru păstrarea mai compactă, într-o formă utilă pentru a fi analizate sau utilizate de aplicații specifice. Agregarea datelor presupune selectarea, filtrarea, combinarea, reorganizarea și manipularea datelor detaliate sau atomice, datelor pentru producerea categoriilor predeterminate și specifice, totalurilor și comparațiilor.

3.1. Date curente detaliate

Interesul cel mai mare îl prezintă datele curente. Datele curente detaliate reflectă cele mai recente evenimente. Aceste date sunt voluminoase și pot ocupa multă memorie pentru păstrare, dar reprezintă, în același timp, cel mai de jos nivel de granularitate. Deseori, ele sunt, pur și simplu, o replică a bazei de date tranzacționale curente, care, mai întâi, este curățată și apoi stocată. Însă nu toate câmpurile păstrate în sistemele tranzacționale pot fi mutate în magazie. Trebuie menționat că, deși sunt referite drept curente, aceste date sunt actualizate o singură dată, în momentul în care sunt trecute în magazia de date.

Datele curente detaliate aproape întotdeauna se păstrează pe disc, care este ușor de accesat, cu toate că administrarea lor este costisitoare și complexă.

3.2. Date vechi detaliate

Datele vechi se păstrează într-o formă masivă de stocare. Ele nu sunt frecvent accesate și sunt păstrate la un nivel de detaliere consistent cu datele curente detaliate. Deși sunt recuperabile în formă detaliată, timpul de acces este mai mare.

3.3. Date slab agregate

Datele slab agregate provin dintr-un nivel jos de detaliere precum nivelul datelor curente detaliate. Acest nivel al magaziei de date aproape întotdeauna se păstrează pe disc.

Experiența arată că agregarea datelor în modul care anticipează aplicațiile îmbunătățește sensibilitatea și utilizarea magaziei de date. Din punctul de vedere al proiectantului, pentru construirea acestui nivel, sunt necesare două decizii: selectarea atributelor și selectarea unităților de timp pentru agregare.

Ambele probleme implică schimbări în care calculele nu trebuie să fie realizate în mod repetat, dar este necesar mult spațiu de păstrare. Este evident că atributele și combinațiile de atribute ce sunt frecvent utilizate în cereri trebuie să fie agregate, în timp ce cele rar utilizate – nu. Odată ce atributele sunt selectate, prima problemă care survine este determinarea frecvenței cu care fiecare atribut va fi agregat.

3.4. Date puternic agregate

Următorul nivel de date, întâlnite într-o magazie de date, sunt datele cu un grad înalt de generalizare. Aceste date sunt compacte și ușor accesibile.

Întotdeauna, unele date, în particular cele solicitate de personalul administrativ superior al unei întreprinderi, trebuie să fie disponibile în formă compactă și trebuie să fie ușor accesibile. Aceste date, de obicei, includ datele ce sunt consultate în mod repetat.

Acest nivel va dispune de capacitatea de menținere a datelor agregate o perioadă lungă de timp, conform tendințelor prestabilite. Odată cu stocarea datelor cu un grad înalt de agregare, timpul de răspuns va fi esențial redus.

3.5. Metadate

Componentul final al structurii datelor, într-o magazie de date, este cel al metadatelor. Metadatele se situează într-o dimensiune diferită de alte date ale

magaziei de date, deoarece conținutul acestora nu este luat direct din mediul operațional.

Metadatele joacă un rol foarte important și sunt utilizate în calitate de:

- directoriu pentru a ajuta analistul în stocarea conținutului magaziei;
- ghid pentru maparea datelor din mediul operațional în magazia de date;
- ghid al algoritmilor utilizați pentru agregarea datelor curente detaliate în date slab agregate și agregarea acestora în date puternic agregate etc.

Metadatele joacă un rol mult mai important într-un mediu data warehouse decât într-un mediu operațional clasic.

Într-o bază de date relațională, metadatele sunt o reprezentare a obiectelor definite în baza de date – în special, ca definiții ale relațiilor, atributelor, bazei de date, viziunilor și altor obiecte. Într-o magazie de date, se referă la ceea ce definește un obiect data warehouse, așa ca relație, atribut, cerere, raport, regulă de afaceri sau algoritm de transformare.

Înțelegerea acestor definiții este esențială pentru toate aspectele procesului de elaborare a unei magazii de date. Gestiunea magaziei de date presupune controlul rigid al tuturor proceselor de la elaborarea programelor, care extrag datele din izvoarele sistemului operațional, la transformarea unei colecții de date într-un obiect al magaziei de date. Magazia de date este utilă, numai dacă aduce un avantaj competitiv, adică, dacă datele transformate pentru umplerea sau stocarea datelor pot fi utilizate pentru a răspunde cererilor de afaceri.

Metadatele reprezintă o hartă sau o schemă a acestor date. Structura metadatelor diferă pentru fiecare proces, deoarece scopul acestora este diferit. Aceasta înseamnă că, în cadrul magaziei de date, există mai multe copii ale metadatelor, care descriu același articol de date. În plus, majoritatea instrumentelor comerciale de administrare a copiilor și acces la date al utilizatorilor finali utilizează propriile versiuni de metadate. Mai exact, instrumentele de administrare a copiilor utilizează metadatele pentru a înțelege regulile de corespondență care trebuie aplicate pentru a transforma datele-sursă într-o formă comună. Instrumentele de acces ale utilizatorilor finali utilizează metadatele pentru a înțelege cum să construiască o interogare. Administrarea metadatelor în cadrul magaziei de date reprezintă o sarcină foarte complexă, care nu trebuie subestimată.

4. ARHITECTURA UNEI MAGAZII DE DATE

Unul din motivele pentru care numărul de magazii de date elaborate crește rapid constă în faptul că

această tehnologie, realmente, este foarte înțeleasă. De fapt, o înmagazinare de date poate reprezenta mai bine structura amplă a unei întreprinderi pentru administrarea datelor din cadrul acesteia. Pentru a înțelege cum interacționează componentele unei

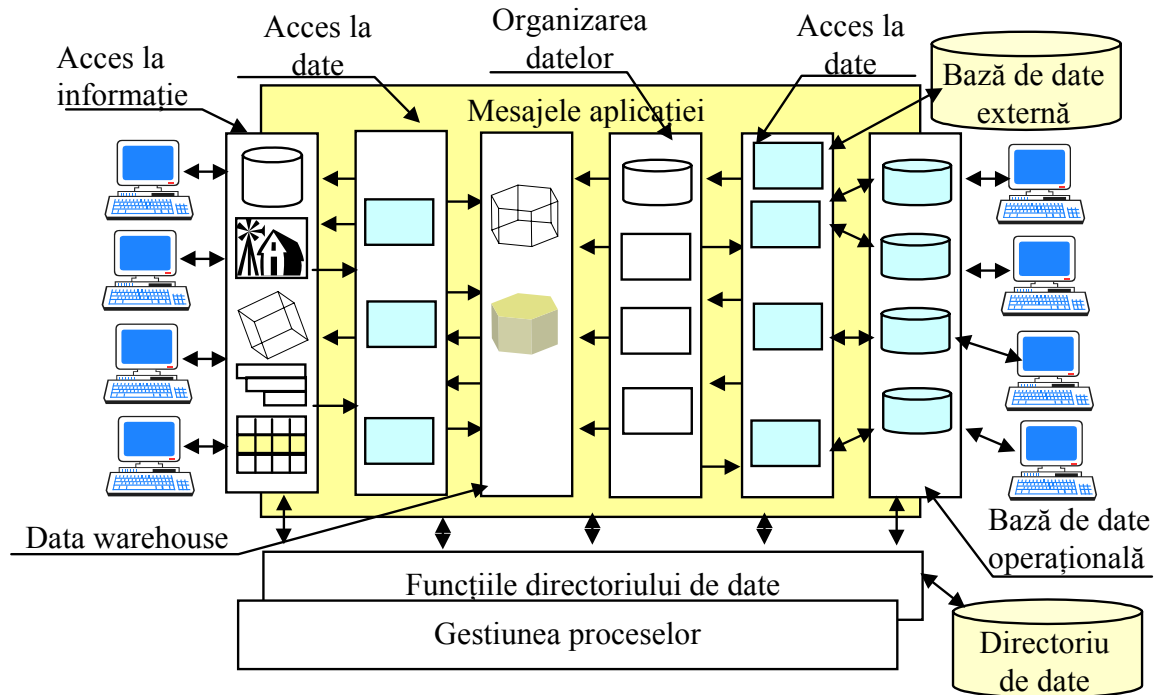


Figura 3. Arhitectura tipică a unui data warehouse.

magazii de date, se trece în revistă arhitectura tipică reprezentată în figura 3.

4.1. Elementele constituente

O arhitectură a unei magazii de date (*Data Warehouse Architecture*) constituie o formă de reprezentare a structurii de date, comunicații, procese și prezentări, care există pentru utilizatorii finali ce dispun de un calculator în cadrul întreprinderii.

Arhitectura este constituită dintr-o mulțime de niveluri interconectate:

- baza de date operațională - nivelul *baza de date externă*;
- nivelul *accesul la informații*;
- nivelul *accesul la date*;
- nivelul *directoriul de date* (metadatele);
- nivelul *gestiunea proceselor*;
- nivelul *mesajele aplicației*;
- nivelul *data warehouse*;
- nivelul *organizarea datelor*.

4.1.1. Baza de date operațională – nivelul *baza de date externe*

Sistemele operaționale prelucrează datele pentru a susține necesitățile operaționale ale întreprinderii. Pentru aceasta, se creează o bază de date operaționale istorice, care asigură o structură de prelucrare eficientă pentru un număr relativ mic de tranzacții comerciale bine definite.

Dar, din cauza tehnicilor limitate ale sistemelor operaționale, bazele de date proiectate pentru suportarea acestor sisteme au dificultăți la accesarea datelor pentru altă gestiune sau scopuri de procesare.

Această dificultate de accesare a datelor operaționale este amplificată și de faptul că multe dintre aceste sisteme au o vechime de 10 – 15 ani. Vârsta unor asemenea sisteme face că tehnologiile de acces la datele disponibile pentru a obține date operaționale sunt, de asemenea, vechi.

Este cert că un scop al magaziei de date este eliberarea datelor ce se păstrează într-o bază de date operațională și combinarea lor cu datele din alt flux de date, în general, extern.

Tot mai mult, organizațiile mari cer date adiționale din baze de date externe. Aceste date includ tendințe demografice, econometrice, achizitive și competitive. Internetul, numit, de asemenea, "*information superhighway*", oferă accesul la multe surse de date în fiecare zi.

4.1.2. Nivelul *accesul la informații*

Nivelul *accesul la informații* al arhitecturii unei magazii de date este nivelul de care utilizatorul final este legat direct. În particular, el reprezintă instrumentele pe care utilizatorul final, de obicei, le utilizează zi de zi. De exemplu, *Excel*, *Lotus 1-2-3*, *Focus*, *Acces*, *SAS* etc.

Acest nivel include, de asemenea, hardware-ul și software-ul implicate în afișarea informațiilor pe ecran și emiterea rapoartelor pentru imprimare, foilor de calcul, graficelor și diagramelor pentru analiză și prezentare. Pe parcursul a două decenii, acest nivel s-a lărgit enorm, în special pentru utilizatorii finali, care au trecut de la PC-urile monoutilizator la PC-urile în rețea.

Actualmente, există instrumente mult mai sofisticate pentru manipularea, analiza și prezentarea datelor. Cu toate acestea, există probleme ce țin de tratarea și convertirea datelor, ca să poată fi recolectate, păstrate și făcute transparente pentru instrumentele și utilizatorii finali. Una din cheile soluționării acestor probleme ar fi elaborarea unui limbaj de date comun ce ar putea fi folosit prin toată întreprinderea.

4.1.3. Nivelul *accesul la date*

Nivelul *accesul la date* al arhitecturii unei magazii de date este implicat, împreună cu nivelul *accesul la informații*, în conversarea cu nivelul operațional. Pe piața mondială, astăzi, limbajul de date comun este SQL. Inițial, SQL a fost elaborat de IBM ca un limbaj de cereri, dar, în ultimii douăzeci ani, a devenit un standard al schimbului de date.

Unul din progresele-cheie ale ultimilor ani a fost elaborarea unei serii de "filtre" de acces la date, precum EDA/SQL. Ele servesc la accesarea aproape la toate sistemele de gestiune ale bazelor de date și sistemele de fișiere de date. Aceste filtre permit instrumentelor de acces accesarea, de asemenea, a datelor păstrate cu sistemele de gestiune a bazelor de date de o vechime de douăzeci de ani.

Nivelul *accesul la date* nu numai conectează SGBD-uri diferite și sisteme de fișiere asupra aceluiași hardware, dar și **celor fabricate și protocoale de rețea**. Una din strategiile-cheie ale

înmagazinării de date este asigurarea utilizatorilor finali cu "un acces la datele universale".

Accesul la date universale semnifică că, din punct de vedere tehnic cel puțin, utilizatorii finali nu țin cont de instrumentele de acces la date sau locație. Ei trebuie să fie capabili să acceseze orice sau toate datele întreprinderii care sunt necesare acestora pentru a-și îndeplini obligațiunile de serviciu.

Nivelul *accesul la date*, atunci, este responsabil de interfața dintre instrumentele de acces la informație și la bazele de date relaționale. În unele cazuri, aceasta este tot de ce are nevoie utilizatorul final. Cu toate acestea, în general, organizațiile elaborează un plan mult mai sofisticat pentru a susține tehnologiile de înmagazinare a datelor.

4.1.4. Nivelul *directorii de date (metadatele)*

Pentru a asigura accesul la datele universale, este absolut necesară menținerea unui directoriu de date sau repository de metadate. Metadatele sunt date despre datele din întreprindere.

Astfel, descrierea unei înregistrări într-un program COBOL sunt metadate. De asemenea, metadate este sentința DIMENSION într-un program Fortran sau sentința de creare în limbajul SQL:

Pentru a avea o magazie de date totalmente funcțională este necesar de a avea o varietate de metadate disponibile, informații asupra viziunilor de date ale utilizatorilor finali și informații asupra bazelor de date operaționale. Ideal, utilizatorii finali trebuie să accedă datele din magazia de date (sau din bazele de date operaționale), fără a avea cunoștințe unde datele rezidă sau forma în care sunt păstrate.

4.1.5. Nivelul *gestiunea proceselor*

Nivelul *gestiunea proceselor* este responsabil de programarea diferitelor sarcini ce trebuie realizate pentru construirea și menținerea magaziei de date și datelor din directoriul de date. Acest nivel poate depinde de multe procese (proceduri) ce trebuie să existe pentru menținerea magaziei de date.

4.1.6. Nivelul *mesajele aplicației*

Nivelul *mesajele aplicației* ține de transportarea datelor în rețeaua întreprinderii. Mesajul aplicației este referit, de asemenea, ca "subproduct", dar poate implica numai protocoale

inconsistente și/sau pot fi codificate în moduri diferite. Toate aceste inconsistențe trebuie rezolvate înainte ca elementele de date să fie acumulate și stocate în magazia de date.

Evident că atât procesele fluxului intern, cât și celorlalte fluxuri necesită crearea metadatelor. Unii cercetători afirmă că chiar există *metafluxul*, adică procesele asociate administrării metadatelor. Metadatele (adică, datele despre date) descriu conținutul magaziei de date. Metadatele constau din definițiile elementelor din magazie, sistemele de date-sursă. În calitate de date, metadatele se integrează și se transformă înainte de a fi stocate.

4.2.2. Fluxul extern

Utilizatorii accesează magazia de date, folosind instrumente ce includ sisteme GUI (*Graphical User Interface*).

Utilizatorilor le pot fi oferite diverse tipuri de instrumente. Acestea pot include software-uri de consultare, generare a rapoartelor, procesare analitică on line, instrumente data/visual mining etc., în funcție de categoria utilizatorilor și cerințele particulare ale acestora. Cu toate acestea, un singur instrument nu poate satisface toate cerințele, prin urmare, este necesară integrarea unei serii de instrumente.

Definiția 4. *Fluxul extern* constă în procesele asociate punerii la dispoziție a datelor pentru utilizatorii finali.

Fluxul extern are loc acolo unde valoarea reală a înmagazinării datelor este percepută de către organizație. Există două activități-cheie implicate în fluxul extern:

- *accesarea*, care se referă la satisfacerea cererilor utilizatorilor finali privind datele de care au nevoie;
- *livrarea*, care e preocupă de livrarea activă a datelor către stațiile de lucru ale utilizatorilor finali.

Magaziile de date, care conțin date agregate, pun la dispoziție un număr distinct de surse de date, pentru a răspunde unei interogări specifice – inclusiv înseși datele detaliate și orice număr de grupuri care satisfac cerințele informaționale ale interogării.

În funcție de aplicație, utilizarea unei magazii de date poate fi extinsă prin capacitatea de accesare a datelor externe. De exemplu, datele accesibile on line prin serviciile computerului sau via Internet, pot fi disponibile utilizatorilor magaziei de date.

4.2.3. Fluxul ascendent

Definiția 5. *Fluxul ascendent* reprezintă procesele asociate adăugării unei valori datelor din magazia de date, prin agregarea, împachetarea și distribuirea datelor.

Activitățile asociate fluxului ascendent cuprind:

- *agregarea datelor*, prin selecția, proiectarea, uniunea și gruparea datelor relaționale în cadrul unor viziuni, care sunt mai convenabile și mai utile pentru utilizatorii finali; agregarea se întinde dincolo de operațiile relaționale simple, pentru a implica o analiză statistică sofisticată, cuprinzând identificarea tendințelor, comasarea și realizarea de eșantioane de date;

- *împachetarea datelor*, prin transformarea datelor detaliate sau agregate în formate mai utile, cum ar fi foile de calcul tabelar, documentele de tip text, diagramele, reprezentările grafice, bazele de date personale și animația;

- *distribuirea datelor* în grupuri adecvate, pentru a mări disponibilitatea și accesibilitatea acestora.

În timp ce se adaugă valoare datelor, este necesar să se acorde atenție și susținerii cerințelor privind performanțele magaziei de date, ca și minimizării continue a costurilor operaționale. În esență aceste cerințe împing proiectarea în direcții opuse. Astfel, administratorul bazei de date trebuie să identifice cel mai adecvat proiect, care îndeplinește toate cerințele, ceea ce, adeseori, necesită anumite compromisuri.

4.2.4. Fluxul descendent

Definiția 6. *Fluxul descendent* reprezintă procesele asociate arhivării și copierii de siguranță a datelor din magazia de date.

Arhivarea datelor vechi, istorice joacă un rol important în menținerea eficacității și performanțelor magaziei de date, prin transferarea datelor mai vechi, cu valori limitate, într-o arhivă de stocare pe benzi magnetice sau discurile optice.

Fluxul descendent de date include procesele de asigurare că starea curentă a magaziei de date poate fi reconstruită, ca urmare a unei pierderi de date sau defecțiuni software sau hardware. Datele arhivate trebuie stocate într-un mod care să permită restabilirea acestora în magazia de date, atunci când este necesar.

Astfel, fluxul de date din figura 4 este normal și prezis într-o magazie de date. Datele intră în magazia de date din mediul operațional (există mici excepții de la această regulă). La intrare în magazia

de date, datele merg spre nivelul de date curente detaliate. Ele rămân aici și se utilizează până apare unul din următoarele evenimente:

- sunt eliminate;
- sunt agregate;
- sunt arhivate.

Apoi odată cu procesul de dezactualizare din magazia de date, datele curente detaliate se mută la nivelul *date vechi*, asociate cu timpul. Procesul de agregare utilizează datele curente detaliate pentru a calcula datele în formă slab sau puternic agregate.

Există mici excepții de la fluxul prezentat. Dar, în general, pentru majoritatea datelor întâlnite într-o magazie, fluxul de date este precum s-a prezentat.

4.2.5. Platforma unei magazii de date

Platforma unei magazii de date este aproape întotdeauna un server de baze de date relaționale. Când se manipulează volume foarte mari de date, poate fi cerută o configurație în bloc de servere UNIX cu multiprocesor simetric (SMP) sau un server cu procesor paralel masiv (MPP) specializat.

Extrasele de date integrate și transformate sunt încărcate în magazia de date. Una din cele mai populare SGBDR-uri disponibile pentru înmagazinări de date asupra platformei UNIX (SMP și MPP), de obicei, este Teradata. Alegerea platformei este importantă. Magazia de date va crește și trebuie să cuprindă cerințele peste 3-5 ani.

Multe organizații cer sau nu alegerea unei platforme din diverse raționamente: “Sistemul X este utilizat de compania concurentă” sau “Sistemul Y este deja disponibil în baza sistemului operațional UNIX, care, deja, este folosit. Una din cele mai mari erori, pe care organizațiile le comit la selectarea platformei, este că ele presupun că sistemul (hardware-ul sau SGBD-ul) se va dezvolta odată cu datele.

4.2.6. Evoluția magaziei de date

Construirea unei magazii de date este o sarcină complexă. Nu se recomandă a înțelege elaborarea unui data warehouse pentru o întreprindere ca o construire a unui proiect ordinar. Mai bine, se recomandă ca cerințele unei serii de faze să fie elaborate și implementate în modele consecutive, care ar constitui un proces de implementare gradual și iterativ.

Nu există nicio organizație care a realizat elaborarea magaziei sale de date, dintr-un singur pas. Multe întreprinderi, însă, au atins acest scop după o elaborare graduală, pas cu pas. Pașii

precedenți evoluează împreună cu materia ce trebuie adăugată.

Datele dintr-o magazie de date nu sunt volatile și reprezintă un depozit de date de o singură lectură (în general). Dar pot fi adăugate elemente noi conform unei baze regulate, pentru ca conținutul să urmeze evoluția datelor în baza de date-sursă atât în conținut, cât și în timp.

Una din problemele menținerii unei magazii de date constă în găsirea metodelor de identificare a datelor noi sau de modificări în bazele de date operaționale. Unele căi de identificare includ inserarea componentei data/timp în înregistrările bazei de date și crearea copiilor de înregistrări actualizate și a copiilor datelor înregistrărilor tranzacțiilor sau ale bazelor de date zilnice.

Aceste elemente de date noi sau modificate sunt extrase, integrate, transformate și adăugate la magazia de date la anumite perioade programate. Când se adaugă noi date, datele vechi sunt eliminate. De exemplu, dacă detaliile unui subiect particular se păstrează 5 ani, când se adaugă datele ultima săptămână, prima săptămână este eliminată.

Bibliografie

1. **Gray Paul, Watson Hugh J.** *Present and Future Directions in Data Warehousing. The Data Base for Advances in Information Systems. Summer 1998, Vol. 29, Nr 3, pp. 83-90.*

2. **Inmon W. H.** *The Data Warehouse and Data Mining. Communications of the ACM. 1996, Vol. 39, Nr 11, pp. 49-57.*

Recomandat spre publicare: 14.11.2012