

## Data Warehouse

Prep. Maria Deac  
Universitatea "Tibiscus" din Timișoara

ABSTRACT. "A data warehouse is a subject-oriented, integrated, time-variant and nonvolatile collection of data in support of management's decision making process." The following article contains a presentation of the Data Warehousing concept and an exemplification of a data warehouse throughout two special programs – DB2 Warehouse Manager and Red Brick Warehouse Manager.

### 1 Prezentare generală

Deoarece în lumea afacerilor de azi doar schimbarea rămâne constantă, dezvoltarea tehnologiei bazelor de date a condus la apariția unor baze de date mari și foarte mari. O astfel de bază de date, în care o organizație păstrează toate datele proprii privind producția, datele financiar contabile, clienții, furnizorii etc. a căpătat denumirea de „*data warehouse*”, adică depozit de date.

Depozitele de informații reprezintă o cerință acută a organizațiilor moderne (fie ele întreprinderi, bănci, administrație, etc) și, totodată, o realitate tehnologică pusă în practică din ce în ce mai frecvent.

Datele din *data warehouse* provin în principal din datele capturate din sistemul operațional, dar mai pot proveni din datele de arhivă (în perioada de constituire a depozitului) precum și din surse externe, cum ar fi baze de date publice.

Câteva exemple posibile: date demografice (obținute în urma unui recensământ), date statistice (furnizate de institute specializate), date de prognoză economică (furnizate de instituții orientate pe studiul pieței), date obținute în urma unor sondaje de opinie, etc. Aceste date pot fi cumpărate, pot fi preluate pe bază de abonament sau pot fi date publice gratuite.

## 2 De ce are nevoie o *data warehouse* ca să funcționeze?

O *data warehouse* trebuie să fie în stare să suporte diferite tipuri de aplicații informaționale. Suportul de procesare a deciziilor este principalul tip de aplicație informațională într-o *data warehouse*, dar folosirea unei *data warehouse* nu trebuie să fie restricționată la un sistem cu suport de decizii. Este posibil ca fiecare aplicație să aibă propriul său set de nevoi în termeni de date, un mod propriu în care datele sunt modelate și felul în care sunt folosite.

O *data warehouse* trebuie să consolideze datele de bază și să furnizeze toate facilitățile pentru a deriva informații din ele după cerințele utilizatorilor. Datele de bază detaliate sunt de primă importanță, dar volumele de date, care tind să fie mari, și utilizatorii au nevoie de obicei de informații derivate din datele de bază. Datele din *data warehouse* trebuie să fie organizate astfel încât să poată fi analizate sau explorate din diferite puncte de vedere.

Din punct de vedere al implementării unui depozit de date există mai mult tendințe: o tendință ar fi implementarea unui sistem distribuit, descentralizat unde datele sunt păstrate în unități independente, „*Independent Data Marts*”, fiecare conținând date relevante pentru un anumit aspect al operațiilor unei instituții. O a doua posibilitate ar fi implementarea unei surse de date unice, centralizate, la care au acces utilizatorii din toate departamentele unei instituții.

Există însă și posibilitatea combinării celor două abordări prin implementarea unei surse de date centralizate la nivelul întregii instituții cu existența unor unități dependente, „*Dependent Data Marts*” care reprezintă subseturi ale datelor din depozitul de date și care au fost selectate și organizate pentru utilizarea specifică anumitor aspecte ale operațiilor unei instituții (resurse umane, contabilitate, etc.). Unitățile dependente *data marts* pot fi implementate logic în cadrul aceluiași depozit de date sau pot fi implementări fizice separate.

Uneori apar greșeli frecvente sau concepții greșite în legătură cu *data marts* și *data warehouse*. O *data marts* sau un „*târg de date*” nu este un depozit de date propriu zis, este un depozit de date cu scopul de a sprijini o aplicație de decizii, care conține date pertinente unui anumit compartiment al unei companii, în timp ce *data warehouse* consolidează toate datele pentru o analiză mai bună, pentru a răspunde oricărei întrebări privind afacerile unei companii.

### 3 Caracteristicile unui depozit de date

Un depozit de date este folosit pentru stocarea datelor financiare non-volatile, tranzacții și evenimente. El constituie o sursă de informații pentru toate suporturile decizionale și aplicații, conținând date integrate consolidate și securizate. În spatele unui *data warehouse* stă un subsistem populat cu un depozit de date obișnuit, o *data warehouse* centrală și un suport de decizii bine definit pentru aplicații informaționale. O sursă pentru un depozit de date poate fi aproape orice sursă relațională sau nerelațională (tabel sau fișier) care este conectată la rețea. Este esențială identificarea tabelelor și fișierelor care vor furniza date în *data warehouse*.

Există mai multe soluții pentru lucrul cu depozitele de date, dintre acestea vom vorbi mai pe larg despre *DB2 Warehouse Manager* din pachetul *DB2 Universal Database* și *Red Brick Warehouse Manager*.

### 4 DB2 Warehouse Manager

*DB2 Warehouse Manager* reunește instrumente pentru a construi, administra și accesa depozitele de date *DB2*. *DB2 Warehouse Manager* simplifică și accelerează crearea prototipurilor de depozite de date, dezvoltarea și implementarea acestora și oferă centrelor de calcul controlul pentru generarea de interogări, analiza costurilor, gestionarea resurselor și urmărirea utilizării. Ajută la satisfacerea cerințelor utilizatorului pentru a găsi, accesa și înțelege informația. Oferă instrumente și tehnici flexibile pentru construcția, administrarea și accesul depozitului de date. *DB2 Warehouse Manager* răspunde celor mai comune cerințe pentru realizarea de rapoarte în întreprinderi de orice mărime. *DB2 Warehouse Manager* se alătură funcțiilor de bază și analitice ale depozitelor de date disponibile în *DB2 Universal Database* oferind:

- scalabilitate sporită a depozitelor de date prin agenți colocați cu baza de date, acești agenți administrează fluxul între sursele depozitelor de date și destinații;
- transformări avansate, utilizând proceduri memorate Java și funcții definite de utilizator, inclusiv filtrarea datelor, pivotarea tabelelor, generarea de chei, etc.;
- un catalog integrat de informații ale afacerii pentru a îndruma utilizatorii către informațiile relevante pe care le pot utiliza în luarea deciziilor;
- guvernarea interogărilor sofisticate și distribuția sarcinilor;

- rapoarte ale interogărilor care satisfac cerințele obișnuite ale majorității întreprinderilor;

## 5 Red Brick

*Red Brick* este un server de baze de date destinat să facă față cerințelor specializate pentru analiza datelor. Furnizează o platformă robustă și scalabilă pentru construcția unor suporturi de aplicație bazată pe decizii. Dă posibilitatea mai multor utilizatori să analizeze mai multe date și să ia decizii mai bune în timp scurt. Procesează seturi mari de date, eficient cu timpi de răspuns foarte mici. Reduce timpul și costul întreținerii unui depozit de date.

Administrarea *data warehouse* se confruntă adesea cu integrarea datelor într-un mediu de baze de date eterogen. Ar putea să existe date stocate în *DB2*, dar ar fi de preferat avantajul capabilității lui *IBM Red Brick Warehouse* de a analiza date. Astfel, ne putem confrunta cu problema transferării datelor din *DB2* în *Red Brick*.

Integrarea acestor două sisteme de baze de date se poate realiza prin *DB2 Warehouse Manager* și aduce avantaje tehnice și beneficii economice. Există anumite metode de configurare care pot să ajute la transferul automat de date foarte mari între *DB2* și *Red Brick*.

Într-un proces de exportare, transformare și încărcare, o bază de date poate fi o sursă de date sau o țintă înspre care merg datele. În *DB Warehouse Manager*, *Red Brick* poate fi configurat să exporte date (ca și sursă) sau să primească date importate (ca și destinație) din alte surse cum ar fi *DB2*.

În afară de a configura o sursă sau o țintă de date *Red Brick*, procesele *Warehouse Manager* care implică *Red Brick* nu sunt diferite de alte surse de date când se folosesc funcții *data warehouse*, cum ar fi transformări, controlul proceselor și programare (schedule).

În cele ce urmează vom dicuta despre două lucruri: În primul rând vom vedea cum se pot exporta date din *Red Brick*. Vorbim despre opțiunile de exportare ale datelor *Red Brick Warehouse* folosind ODBC sau folosind comenzi SQL de export din *Red Brick*.

În al doilea rând vom vedea cum se importă datele în *Red Brick*, adică despre felul cum se importă datele în *Red Brick Warehouse* folosind ODBC sau folosind “încărcătorul *Red Brick*” numit Table Management Utility (TMU) printr-un program definit de utilizator (User Defined Program (UDP)).

## 6 Exportul datelor în Red Brick

Se pot exporta date din *Red Brick Warehouse* în două moduri:

- **Folosind ODBC** - O bază de date *Red Brick* poate fi definită ca și o sursă generică ODBC în *DB2 Warehouse Manager* folosind următorii pași:
  - mai întâi trebuie definit un nume al sursei de date ODBC (Data Source Name DSN) pentru a accesa baza de date *Red Brick*;
  - apoi în *DB2 Data Warehouse Center* se execută pașii: Warehouse Source -> Define ->Generic ODBC pentru a termina pașii de configurare standard;
  - după aceea se poate folosi un pas SQL în *Warehouse Manager* pentru a primi date din orice tabel *Red Brick*, prin interfața ODBC.
- **Folosind exportul SQL** - Red Brick :
  - se execută pașii descriși anterior (pentru ODBC);
  - sintaxa pentru exportul detaliat poate fi găsită în ghidul de referință Red Brick SQL :

**export to output\_file ddl\_file tmu\_file format\_type**

Comanda de export este invocată prin ODBC la fel ca și orice interogare normală, dar scrie direct rezultatul execuției interogării într-un fișier, director sau într-un pipe, trecând peste straturile ODBC. Formatul de fișiere suportate poate fi de lungime fixă și binară. Suportul XML a fost adăugat în *Red Brick* versiunea 2. Toate fișierele de ieșire rezultate în urma exportului SQL trebuie să fie accesibile de pe stația locală.

Atunci când se specifică un format extern, fiecare coloană din rezultat va avea o lungime fixă.

Din cele două metode de export, metoda SQL este de preferat din cauza performanței mai ridicate. Este de reținut faptul că fișierul care rezultă dintr-un export SQL trebuie să fie accesibil de pe stația locală. Dacă utilizatorul trebuie să mute datele între stații, el poate folosi un pas ajutor ftp pentru a trimite fișierele către o stație aflată la distanță.

Metoda ODBC este mai lentă decât exportul SQL, dar poate muta date între stații fără a avea nevoie de pași în plus.

## 7 Importul datelor în Red Brick

La fel ca și exportul datelor se poate opta între două metode de import a datelor.

- **Folosind ODBC** – O bază de date *Red Brick* poate fi definită ca o țintă de date ODBC.
- **Folosind un program TMU- UDP** – utilizatorii pot configura un program UDP în *Warehouse Manager* pentru a folosi *Red Brick* TMU. Pentru a seta parametrii UDP, un utilizator specifică calea până la: un executabil TMU, un fișier de control TMU, o bază de date țintă, un username și o parolă. După aceasta TMU UDP va fi disponibil pe panoul de instrumente al ferestrei modelelor de procese.

Din cele două metode se recomandă TMU UDP pentru că lucrează cu volume mari de date. Această metodă are nevoie ca fișierele de intrare și fișierele TMU să fie accesibile de pe stația locală. Dacă accesul este o problemă se poate folosi clientul TMU din *Red Brick 6.2* sau un pas adițional ftp.

Metoda ODBC poate muta date între stații, dar pentru că este mult mai lentă este folosită pentru lucrul cu volume mai mici de date.

În concluzie, din cauza performanței ridicate, se recomandă folosirea comenzilor de export SQL și a programului TMU UDP pentru exportul/importul în/din baze de date mari. Amândouă metode oferă instrumente pentru întreținerea magaziei de date .

## Bibliografie

- [Inm03] **William H. Inmon**, *Building the Data Warehouse*, third Edition, 2003
- [QiJ03] **Qi Jin**, *IBM Red Brick @ Warehouse Development*, IBM Silicon Valley lab, 2003
- [Sar98] **Mircea Sârbu**, *În căutarea definiției*, Computer Press Agora SRL, 1998