

# Business Intelligence

## Cursul 3



Prof. Ramona Bologna,  
ASE Bucuresti

# Agenda

---

## 1. OLAP

- Cerinte functionale OLAP
- Arhitecturi OLAP: ROLAP, MOLAP, desktop OLAP si Hybrid OLAP

## 2. Data mining

## 3. Integrarea datelor

- BI si ERP;
- Descrierea unui sistem ERP (SAP ) integrat cu software BI

# 3. OLAP



- Cerinte functionale OLAP
  - Regulile lui Codd
  - Regulile FASMI
- Arhitecturi OLAP: ROLAP, MOLAP, desktop OLAP si Hybrid OLAP

# Ce este OLAP

---

- De instrumentele de interfata si structura BD suporta analiza multidimensionala, acces instantaneu si manipulare usoara => **online analytical processing**
- **Codd**, parintele acestui termen a evidentiat diferentele OLTP-OLAP- 1993 - **criterii generale** pentru BD OLAP.
- **ANALIZA MULTIDIMENSIONALA**
  - Aplicarea de **formule și modele** asupra dimensiunilor și ierarhiilor;
  - Vizualizarea datelor prin **mai multe filtre sau dimensiuni** in acelasi timp
  - **Analiza de trend** pe perioade diferite de timp;
  - Analiza în adancime (**drill-down**);
  - **Extragerea unui subset** de date pentru vizualizare(slice);
  - **Rotații** în cadrul dimensiunilor;

# OLAP si DW

---

- Sistemele OLAP și DW - sisteme suport de decizie orientate pe date și sunt similare.
- **DW** pune accentul pe procesele ce asigură consistența, corectitudinea și valabilitatea datelor la utilizatori,
- **sistemele OLAP** pun accentul pe *cerințele analitice* și *procesele de modelare și calcul* necesare.

# Cerinte functionale OLAP - Codd

## □ Caracteristici de bază

- 1: O viziune conceptuală multidimensională
- 2: Manipularea intuitivă a datelor
- 3: Accesibilitate
- 4: Surse de date variate
- 5: Modele de analiză OLAP
- 6: Arhitectura client/server
- 7: Transparență
- 8: Suport multiutilizator

## Caracteristici speciale

- 9: Denormalizarea datelor
- 10: Stocarea rezultatelor generate de instrumentul OLAP
- 11: Manipularea valorilor lipsă
- 12: Modul de tratare a valorilor lipsă

## Modul de prezentare a datelor

- 13: Flexibilitatea rapoartelor
- 14: Performanța raportării
- 15: Ajustarea automată a nivelului fizic

## Controlul dimensiunilor

- 16: Dimensionalitate generică
- 17: Dimensiuni și niveluri de agregare nelimitate
- 18: Operații între dimensiuni nerestrictive

# Regulile FASMI (1995 Nigel Pendse )

---

- ❑ **Fast Analysis of Shared Multidimensional Information**
- ❑ **FAST** - cat mai multe raspunsuri utilizatorilor intr-un termen mai scurt de 5 sec
- ❑ **ANALYSIS** - orice logica de afaceri si analiza statistica relevanta pentru aplicatie sau utilizator, suficient de simplu pentru utilizatorul final
- ❑ **SHARED** - toate cerintele de securitate pentru confidentialitate, dar si blocarea actualizarilor concomitente, daca este necesar accesul multiplu la scriere
- ❑ **MULTIDIMENSIONAL** - viziune **conceptuala multidimensionala** asupra datelor, inclusiv suport complet pentru ierarhii si ierarhii multiple
- ❑ **INFORMATION** reprezinta toate datele si informatiile derivate necesare oriunde se afla si in orice masura este relevanta pentru aplicatie.

# a. Arhitecturi OLAP - ROLAP

---

- ❑ **Relational OnLine Analytic Processing**
- ❑ tehnologia relațională, adaptată și extinsă
- ❑ agregările sunt stocate în cadrul BD relaționale sursă
- ❑ cea mai lentă soluție, ex: **DSS Server/Microstrategy**
- ❑ **Avantaje:**
  - se integrează cu tehnologia și standardele existente;
  - actualizarea sistemelor MOLAP este dificilă;
  - ROLAP sunt adecvate pentru a stoca volume mari de date, prin utilizarea **procesării paralele** și a **tehnologiilor de partiționare**;
  - ROLAP sunt recomandate pentru aplicațiile cu volatilitate ridicată a datelor (antecalcul agregari)



# ROLAP atunci cand:

---

- a) Volumul de date este **prea mare** pentru a fi duplicat.
- b) Datele sursă **se modifică frecvent** și este mai bine de a citi în timp real decât din copii;
- c) Se dorește **integrarea cu alte sisteme** informatice relaționale existente;
- d) Firma are o **politică de neduplicare** a datelor, pentru securitate sau alte motive, chiar dacă aceasta conduce la aplicații mai puțin eficiente

## b. Arhitecturi OLAP - MOLAP

---

- ❑ **Multidimensional OnLine Analytic Processing**
- ❑ stocarea datelor în formă **multidimensională**, folosind **structuri de date vector (tehnica matricilor rare)**
- ❑ atât datele sursă, cât și agregările sunt stocate în format multidimensional
- ❑ indexare rapidă a datelor preagregate
- ❑ opțiunea cea mai *rapidă pentru consultare*
- ❑ necesită cel mai *mult spațiu* de disc
- ❑ stocarea fizică a datelor multidimensionale, precum și fenomenul de împrăștiere sunt preocupări majore
- ❑ Ex: **Oracle Essbase**

# Avantaje MOLAP

---

- tabelele nu sunt potrivite pentru date multidimensionale;
- **matricile multidimensionale** permit stocarea eficientă a datelor multidimensionale;
- **limbajul SQL** nu este corespunzător pentru operații OLAP

## c. Arhitecturi OLAP - HOLAP

---

- ❑ **Hybrid OnLine Analytic Processing**
- ❑ combinație a primelor două modele
- ❑ Arhitecturi HOLAP
  - agregările - stocate in structură multidimensională, nivelul celulelor de bază în formă relațională
  - cele mai recente felii de date stocate in MOLAP si restul in ROLAP
- ❑ oferă performanțele MOLAP atunci când este nevoie de preluarea datelor din tabele
- ❑ Ex: [Microsoft SQL Server OLAP Services](#)

# Caracteristici HOLAP

---

- transparența locației și a accesului
- transparența fragmentării
- transparența performanței
- un model de date comun
- alocarea optimă în sistemele de stocare

## Solutii OLAP

	Essbase	SSAS	Mondrian	Oracle OLAP	SAS OLAP	SAP Netweaver BW	Microstrategy Intelligence Server	Cognos TM1
<b>Firma</b>	Oracle	Microsoft	Pentaho	Oracle	SAS Institute	SAP	Microstrategy	IBM
<b>Arhitectura</b>	MOLAP ROLAP HOLAP	MOLAP ROLAP HOLAP Power Pivot pt excel	ROLAP	MOLAP ROLAP HOLAP	MOLAP ROLAP HOLAP	MOLAP ROLAP	MOLAP ROLAP HOLAP	TM1- MOLAP
<b>Lb MDX</b>	DA	DA	DA	NU	DA	DA	DA	DA
<b>Lb SQL</b>	NU	DA	NU	DA	NU	NU	DA	NU
<b>Standard XML for analysis</b>	DA	DA	DA	NU	DA	DA	DA	DA
<b>Masuri semiaditive</b>	DA	DA	NU	DA	DA	DA	DA	DA
<b>SO</b>	Win Linux Unix	Win	Win Linux Unix	Win Linux Unix	Win Linux Unix	Win Linux Unix	Win Linux Unix	Win Linux Unix

# Operatori OLAP: Slice, dice, pivot

---

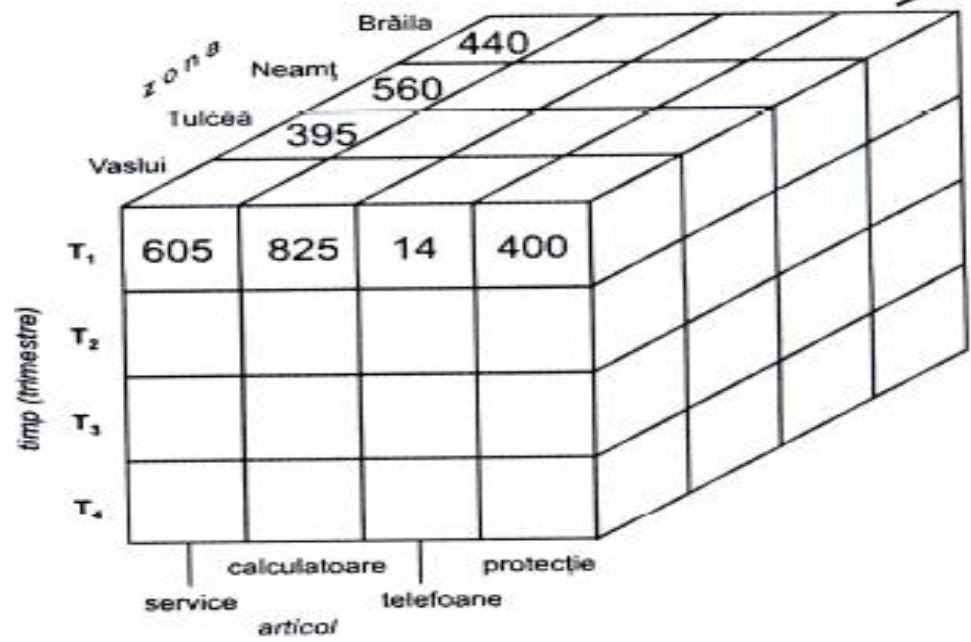
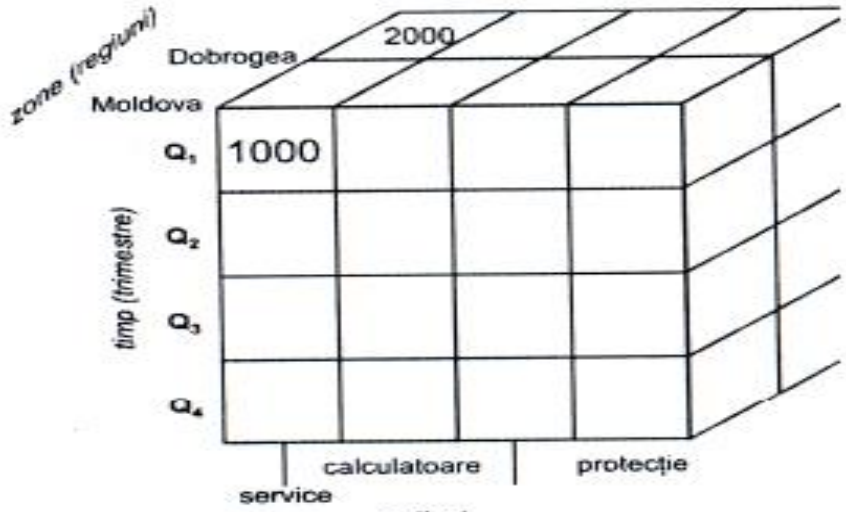
- ❑ **Secționarea (slice)** –selectarea unui membru a unei dimensiuni, crearea unei serii de intersecții cu alte dimensiuni al acelui membru pentru secțiunea respectivă.
- ❑ **Decupare (dice)** – crearea unui subcub al cubului de date prin selectarea unor dimensiuni și a intervalelor de valori pentru acestea. De exemplu, vânzări după luna, după regiune, după client. Acest “după” ne indică cum putem realiza **rotirea (dicing)** datelor.
- ❑ **Pivotarea și imbricarea** dimensiunilor.  
**Pivotarea** presupune înlocuirea între ele a dimensiunilor în cadrul unei vizualizări, trecerea de pe linii pe coloane și invers.

# Ierarhii si navigare

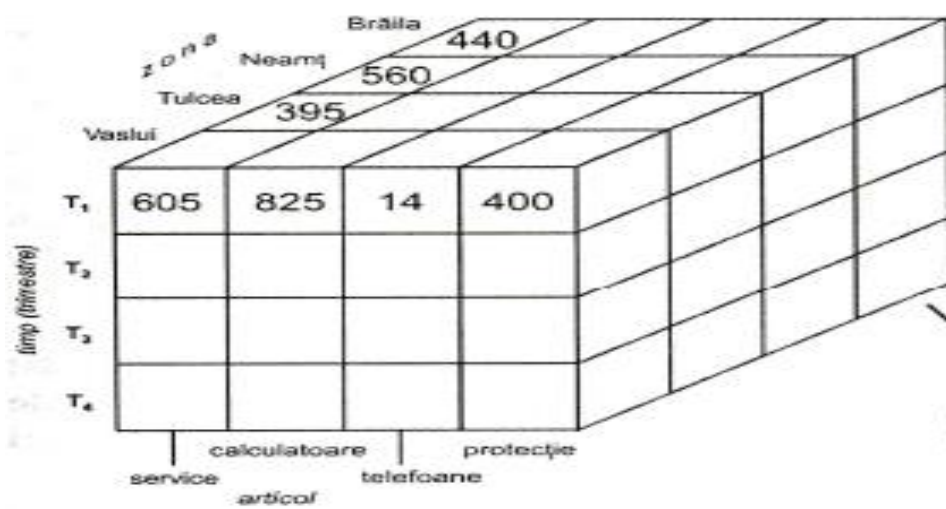
---

- ❑ Sistemele OLAP permit agregarea sub formă de ierarhii care realizează agregarea detaliilor la nivele din ce în ce mai înalte.
- ❑ De exemplu, datele lunare pot fi agregate (**roll-up**) la nivel de trimestru sau la nivel de an,
- ❑ Plecarea din vârful ierarhiei dimensiunilor și detalierea (**drill-down**)
- ❑ Ia nastere astfel un nou proces, denumit «**analiză ad-hoc**», care permite :
  - Răspunsul la diferitele întrebări ale managerilor în doar câteva minute de navigare în date
  - Formatarea rapoartelor prin pivotarea și imbricarea dimensiunilor
  - Învățarea rapidă a utilizării unui astfel de sistem de către orice persoană, mai ales managerii

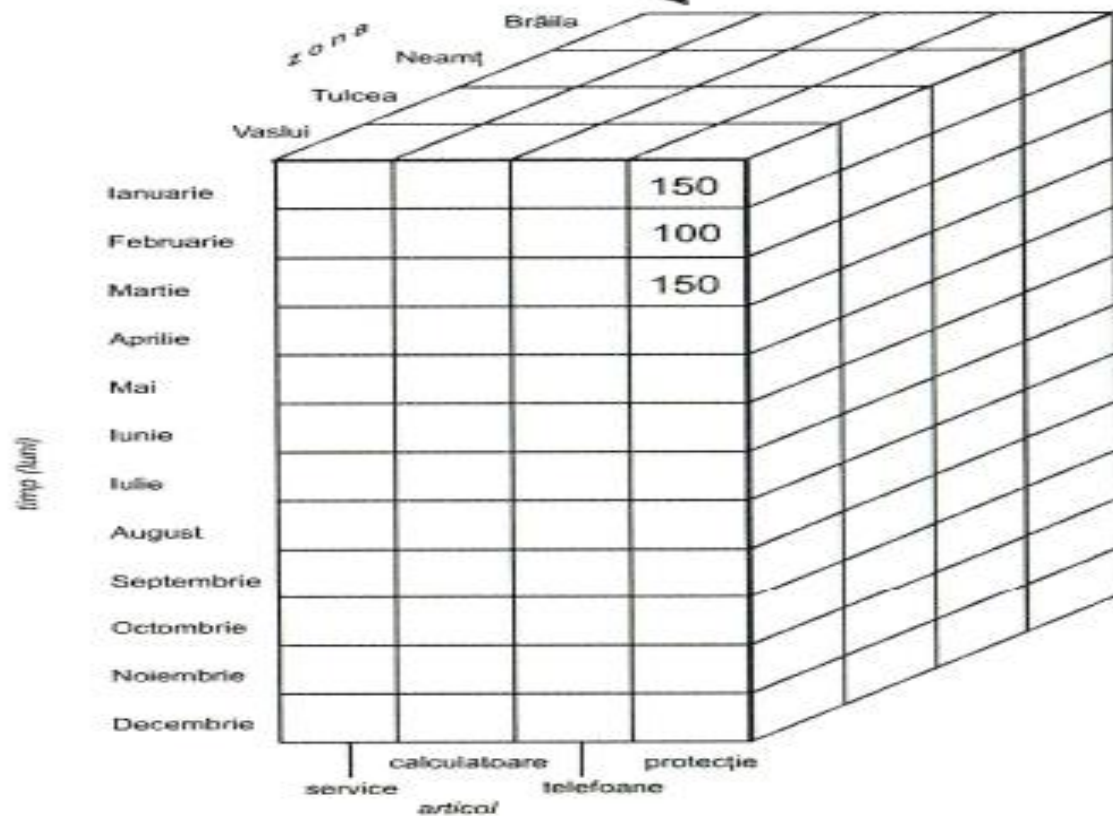


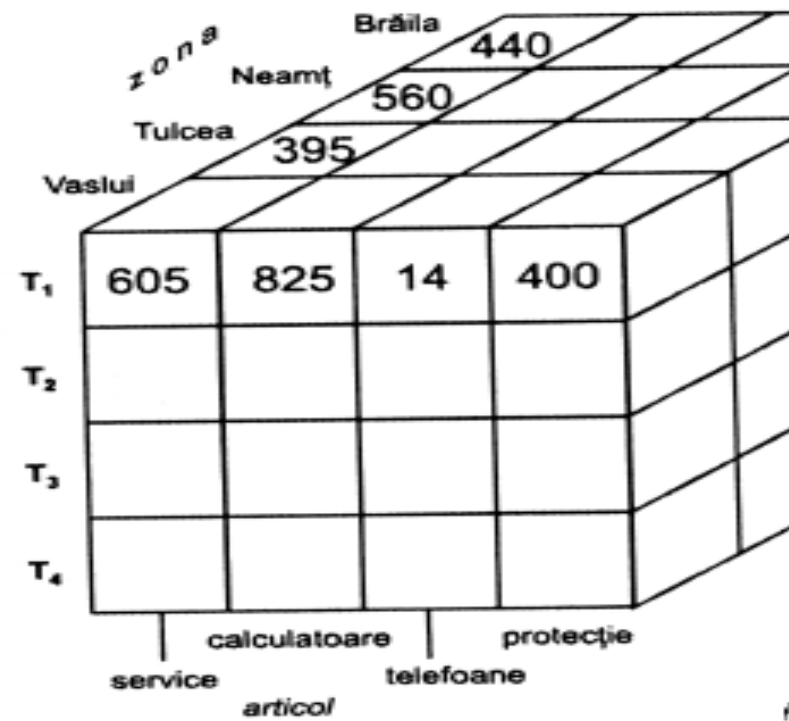


roll-up  
on zona  
(from localitate  
to regiune)



drill-down on timp (from trimestru to luna)

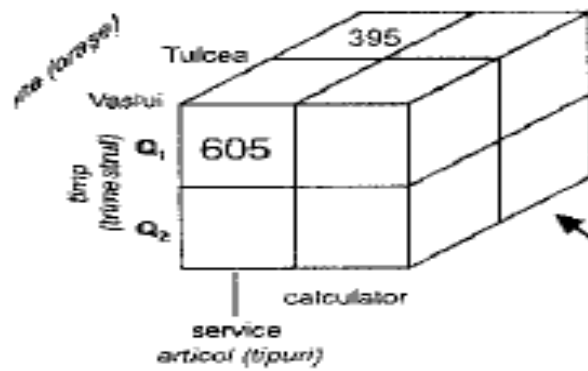




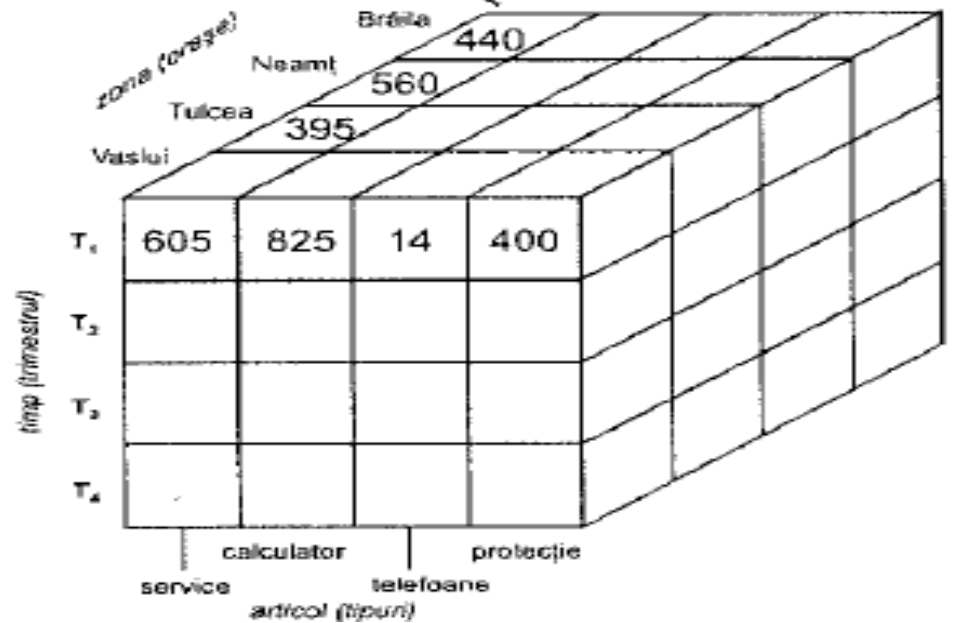
slice for timp="T<sub>1</sub>"

zona

Brăila				
Neamț				
Tulcea				
Vaslui	605	825	14	400
	service	calculatoare	telefoane	protecție
		articol		

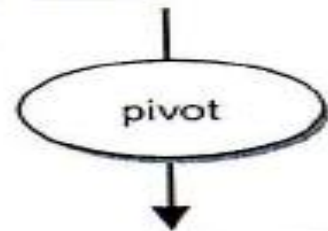


dice for  
 (zona="Tulcea" or "Vaslui")  
 and (time="Q1" or "Q2") and  
 (articol="service" or "calculator")



zona	Brăila				
	Neamț				
	Tulcea				
	Vaslui	605	825	14	400

| calculatoare | protecție  
 service      telefoane  
 articol



articol	service				605
	calculatoare				825
	telefoane				14
	protecție				400
		Brăila	Neamț	Tulcea	Vaslui
		zona			

## 2. Data Warehouse si Data mining

---

- Trei tipuri de aplicatii de DW
  - **Procesarea informatiilor**
    - Interogari, analize statistice de baza, raportari folosind tabele, grafice, figuri
  - **Procesare analitica**
    - Analiza multidimensionala a datelor DW
    - Operatii OLAP de baza, navigare prin date, pivotari, rotatii, sectionari
  - **Data mining**
    - Descoperire de cunostinte din modele ascunse
    - Asocieri, construire de modele analitice, realizare de clasificari si predictii, si prezentarea rezultatelor cu instrumente de vizualizare
    - OLAM –Online Analytical Data Mining

# Despre data mining

---

- ❑ Converteste **datele** in **cunostinte valoroase** care pot fi folosite ca suport pentru decizii
- ❑ Este o *colectie de metodologii, tehnici si algoritmi de analiza a datelor* pentru descoperirea de **modele noi** in date
- ❑ Este folosit pentru **seturi mari de date**
- ❑ Procesul este **automatizat**, nu e necesara interventia umana
- ❑ ***Data mining*** si ***Knowledge Discovery in Databases (KDD)*** sunt considerate de unii autori ca reprezentand acelasi lucru. Altii considera data mining-ul ca fiind *pasul de analiza* in procesul KDD, dupa curatarea si transformarea datelor si inainte de vizualizare/ evaluarea rezultatelor

# Despre data mining

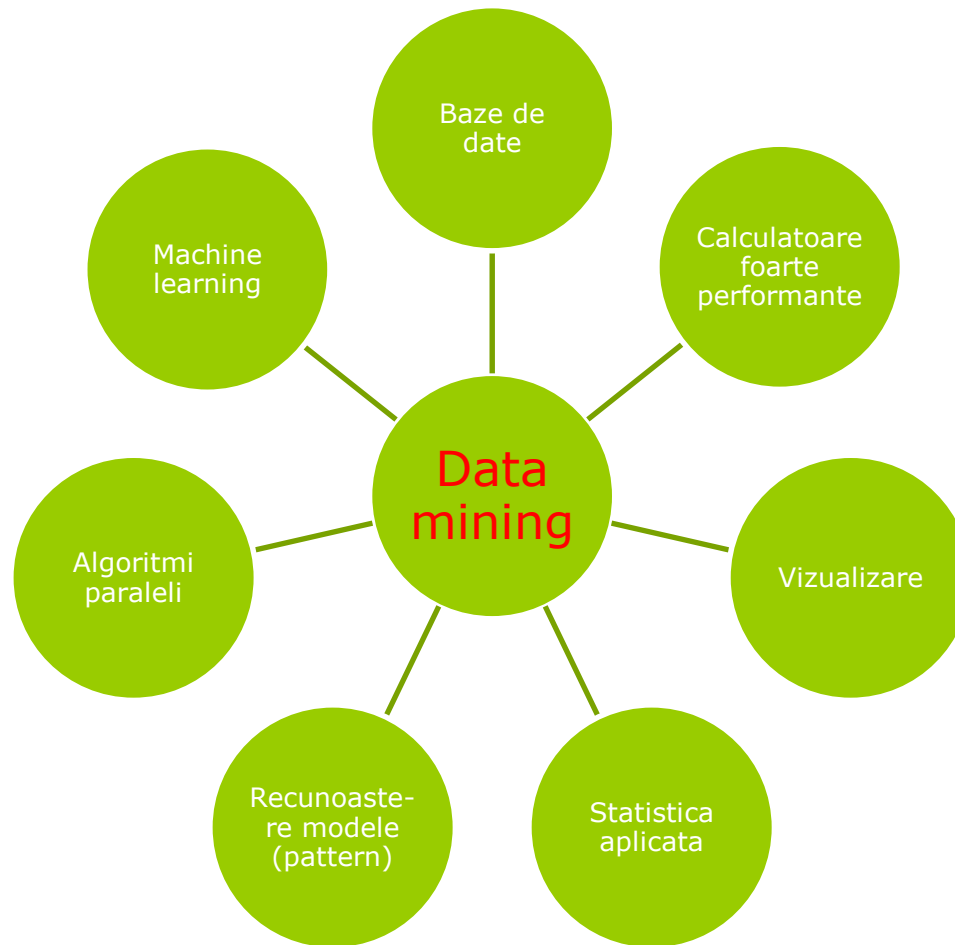
---

- ❑ Modelele trebuie sa fie valide, utile si inteligibile
- ❑ Implica **metode** care sunt la intersectia intre inteligenta artificiala, invatare automata (machine learning), statistica si sisteme de baze de date.
- ❑ Cele mai valoroase **rezultate** obtinute prin DM sunt: clusterizarea, clasificarea, estimarea, predictia si gasirea lucrurilor care apar impreuna.
- ❑ Principalele instrumente de DM includ:
  - Arbori de decizie,
  - Retele neuronale
  - Instrumente de vizualizare
  - Algoritmi genetici
  - Logica fuzzy
  - Metode statistice clasice.



# De unde provine?

---



# Povesti de succes

---

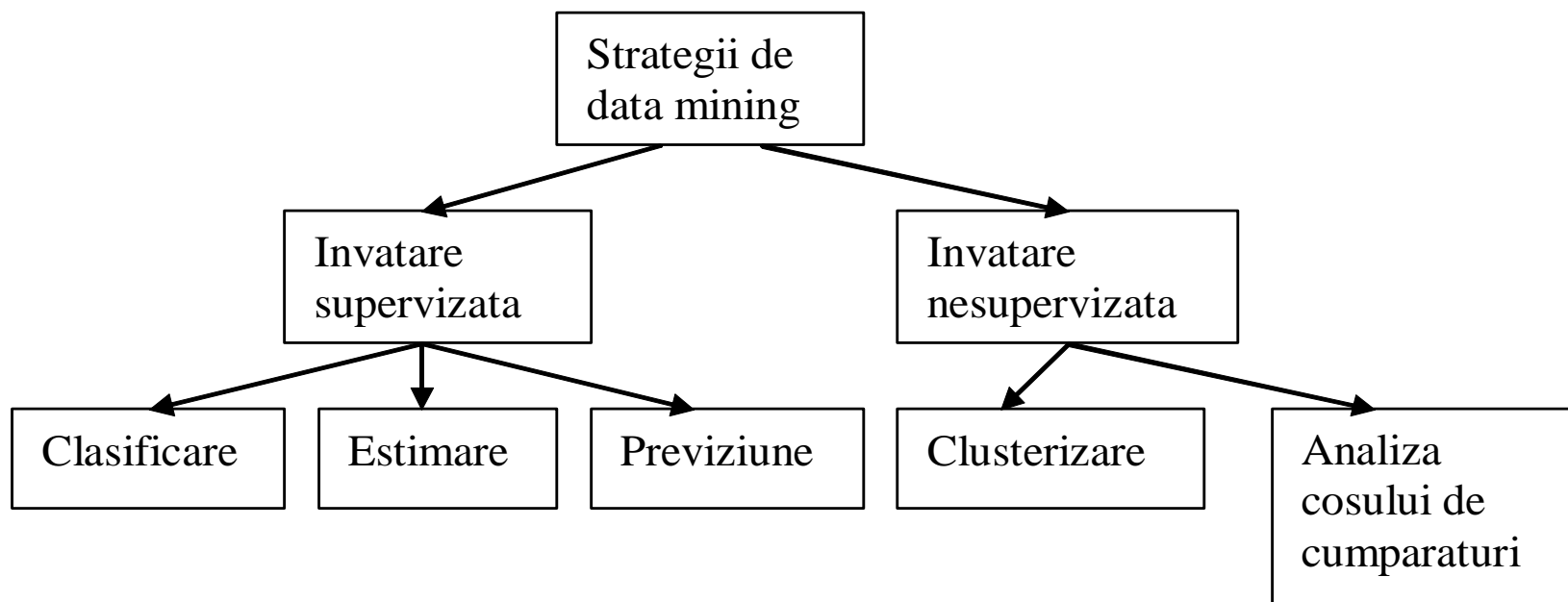
- ❑ **Arbori de decizie** construiti pe baza istoricului **imprumuturilor** acordate de banci pentru a decide daca sa acorde un imprumut
- ❑ **Modele ale comportamentului de calatori** utilizate pentru a gestiona vanzarea cu discount a biletelor de avion, camerelor de hotel, etc.
- ❑ **"Scutece si bere"** Observatia ca acei clienti care cumpara scutece au probabilitate mai mare decat media de a cumpara si bere a permis supermarket-urilor sa amplaseze Berea langa scutece. Si chips-urile langa. Au crescut vanzarile la toate trei.
- ❑ Skycat si Sloan Sky Survey: **clusterizarea corpurilor ceresti** in functie de nivelul lor de radiatii in diferite benzi, a permis astronomilor sa diferentieze intre galaxy, stele in formare si alte tipuri de obiecte ceresti.
- ❑ **Compararea genotipului** unor persoane care au/ nu au anumite afectiuni au permis descoperirea unor seturi de gene responsabile pentru cele mai multe din cazurile de diabet. Acest gen de DM va deveni tot mai util, odata ce genomul uman este construit

# Data mining

---

## □ Pasi:

- i: Culegerea datelor
- ii: Preprocesarea (pregatirea) datelor de analizat.
- iii: Analiza datelor sau aplicarea unui algoritm/metode de DM
  - Invatare supervizata
  - Invatare nesupervizata
- iv: Vizualizarea si interpretarea rezultatelor algoritmului
- v: Aplicarea rezultatelor obtinute la noi probleme.



# Pasii procesului de data mining

---

1. **Culegerea datelor:** colectarea datelor din baze de date sau prin cautari pe Web
2. **Preprocesarea datelor**
  - a) **Curatarea datelor:** inlocuirea (sau stergerea) valorilor lipsa, eliminarea sau doar identificarea valorilor extreme, eliminarea zgomotelor din date, inlaturarea inconsistentelor.
  - b) **Integrarea datelor:** datele sunt preluate din surse multiple, cu tipuri de date si structure diferite, sunt integrate si se elimina duplicatele si inconsistentele
  - c) **Transformarea datelor:** normalizarea (sau standardizarea), sumarizari, generalizari, construirea de noi attribute, etc.
  - d) **Reducerea datelor** (sau extragerea caracteristicilor): doar attributele relevante sunt selectate pentru procesare ulterioara
  - e) **Discretizarea:** deoarece unii algoritmi lucreaza doar cu valori discrete, valorile atributelor continue trebuie inlocuite cu valori discrete dintr-o lista predefinita

# Pasii procesului de data mining

---

- 3. Analiza datelor:** se aplica algoritmii de DM si se realizeaza extragerea si descoperirea de modele.
- 4. Vizualizarea:** deoarece DM extrage proprietati si informatii ascunse din date, pentru a intelege si evalua rezultatele e necesara vizualizarea lor
- 5. Evaluarea rezultatelor:** nu toate rezultatele obtinute prin DM sunt informatii valoroase. Pot rezulta adevaruri statistice sau informatii care nu sunt utile in activitatea analizata. Expertii sunt cei care vor evalua rezultatele.

# Metode de data mining

---

- **Metode predictive** – utilizeaza niste variabile pentru a prezice valoarea altor variabile.
  - **Clasificarea** – se bazeaza pe date cunoscute, etichetate si algoritmi de clasificare construiesc modele pentru a clasifica date noi
  - **Regresia**
  - **Detectarea deviatilor** (ex: fraude, intruziuni)
- **Metode descriptive**: algoritmi gasesc modele care descriu structura interna a setului de date.
  - **Clusterizarea** – identifica grupuri de obiecte similare (CLUSTERE) din setul de date, dar si posibile obiecte izolate, valorile extreme.
  - Descoperirea **regulilor de asociere**
  - Descoperirea **pattern-urilor secventiale**

# Intrebare pentru voi!

---

- Care sunt diferentele intre Data mining si Machine learning?
- Exemple

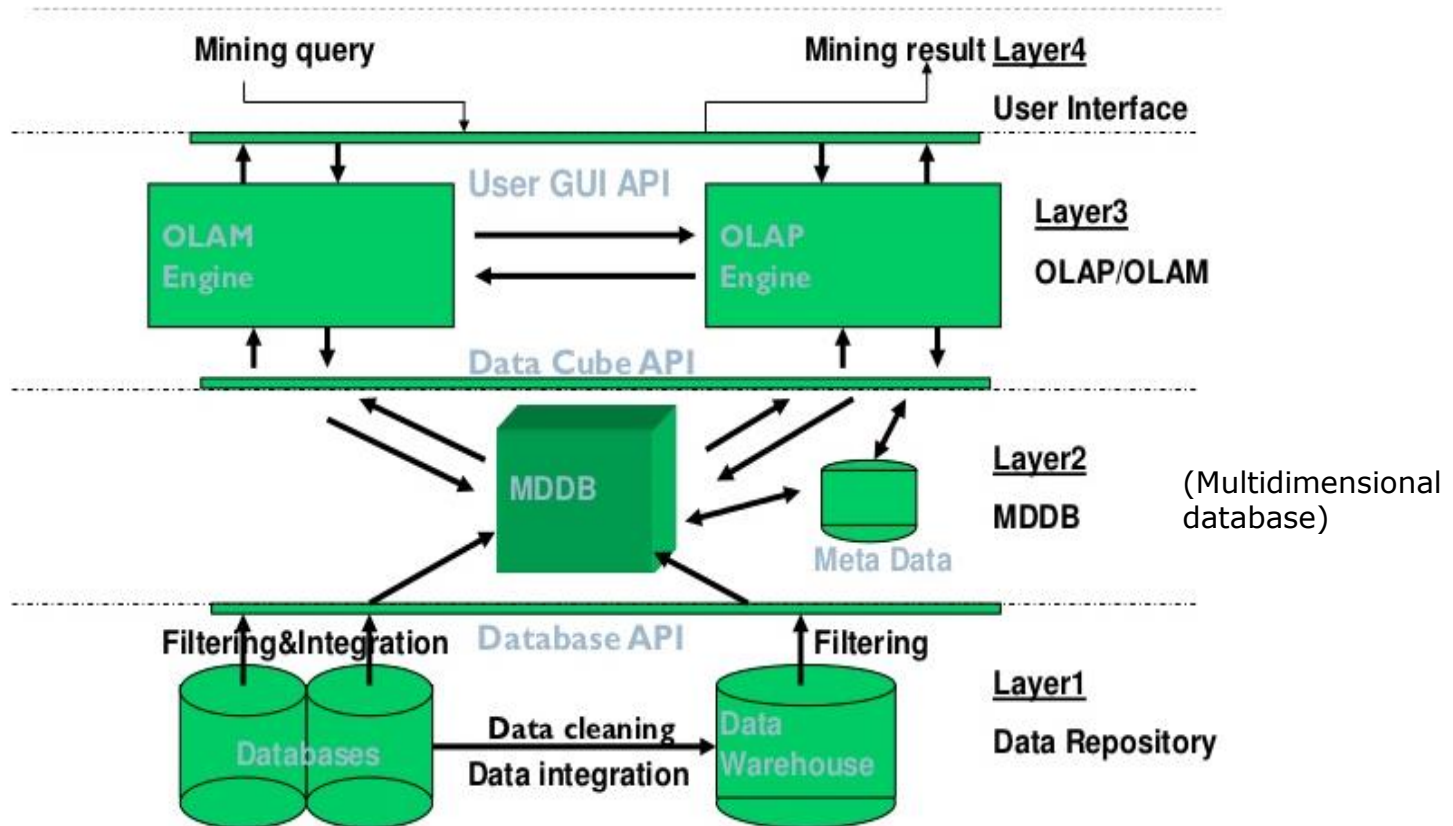
# DM in DW

---

- **Volume f. mari de date** – milioane de inregistrari, mii de attribute
- Se realizeaza **procesul ETL** si se incarca si gestioneaza datele in sistem multidimensional
- Se ofera acces utilizatorilor de business care isi vor realiza **analizele dorite** prin aplicatii software specifice
- **Rezultatele** sunt prezentate sub forma de tabele sau grafice



# Arhitectura sistem OLAM



# Aplicatii DM

---

- ❑ **AT&T** utilizeaza o aplicatie de data mining pentru identificarea apelurilor internationale frauduloase;
- ❑ sistemul american **FAIS** (Financial Crimes Enforcement Network AI System) utilizeaza data mining pentru identificarea activitatilor de spalare a banilor in cadrul tranzactiilor foarte mari de bani;
- ❑ **Banca Americii** utilizeaza data mining pentru identificarea clientilor care utilizeaza anumite produse ale bancii si care sunt produsele preferate ale clientilor, in scopul crearii de mixuri de produse care sa satisfaca exigentele clientilor.
- ❑ **US West Communications**, furnizor de servicii de comunicatii cu peste 25 milioane de clienti, utilizeaza data mining pentru a determina tendintele si nevoile clientilor pe baza unor parametri de tipul: dimensiunea familiei, varsta medie a membrilor familiei si adresa de rezidenta.
- ❑ **Twentieth Century Fox** analizeaza incasarile de box-office pentru a identifica care actori, filme si scenarii vor fi apreciate in diverse arii de marketing.

# 3. Integrarea datelor

---

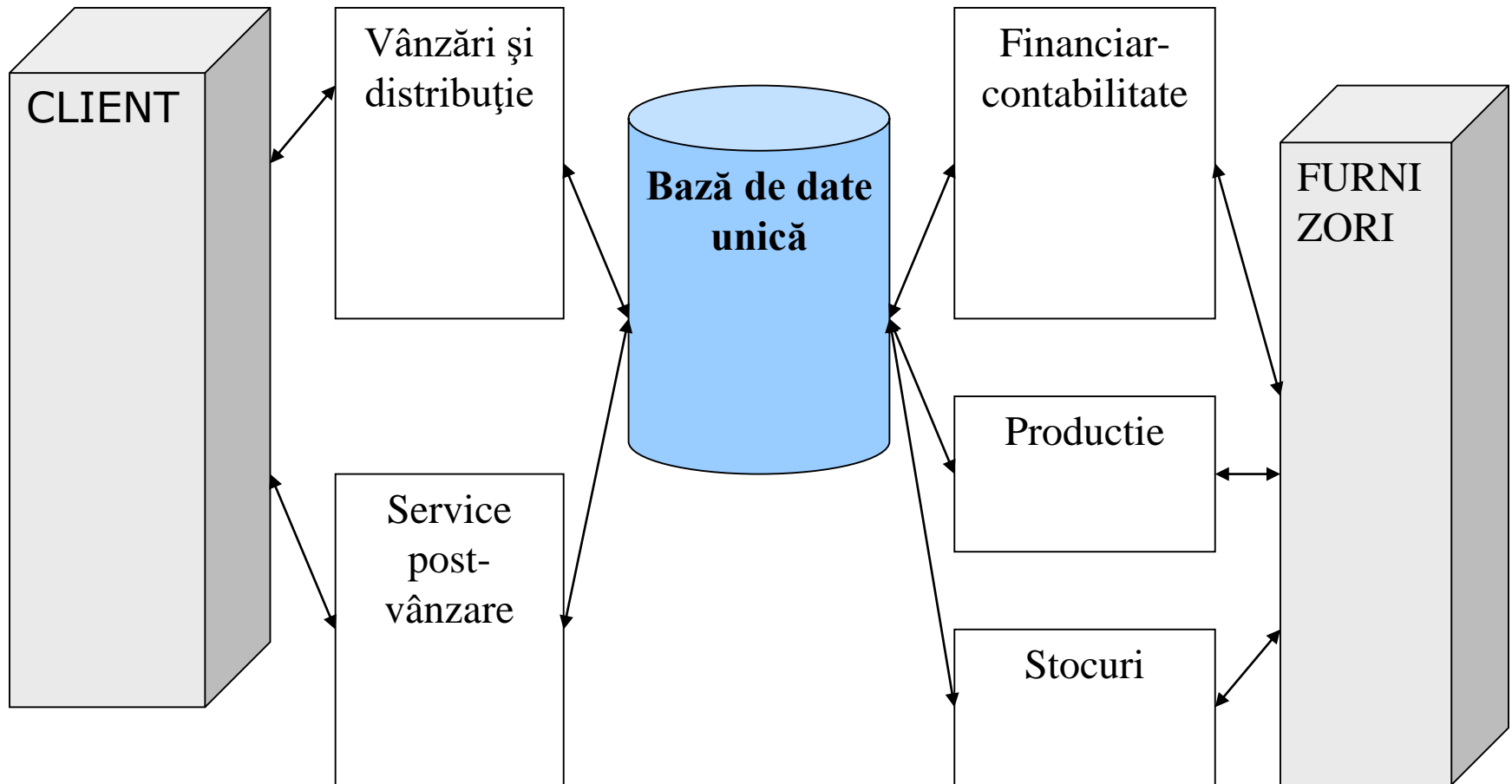
- BI si ERP;
- Descrierea unui sistem ERP (SAP ) integrat cu software BI

# Ce este un sistem ERP

---

- “un pachet care promite **integrarea completă** a tuturor informațiilor din cadrul unei organizații” [Davenport]
- “infrastructură software, **multimodulara** ce oferă suport de gestiune și coordonare a diferitelor structuri și procese din companie, în vederea realizării obiectivelor de afaceri” [Fotache]
- Oferă **accesabilitate, vizibilitate și consistența informațională** în întreaga organizație
- Dezvoltare cu instrumente CASE

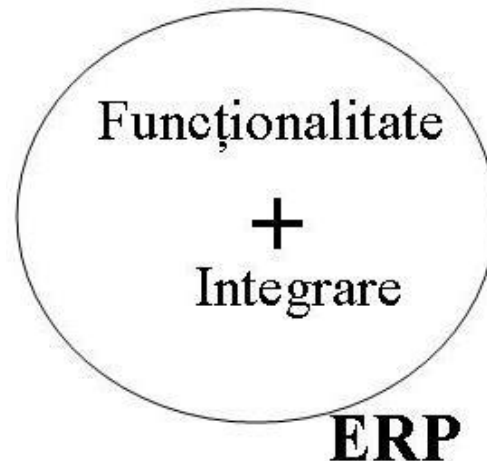
# Arhitectura client-server



# Proprietati fundamentale

---

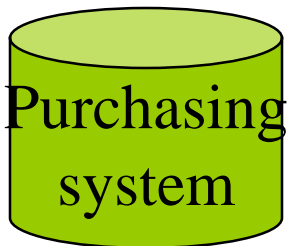
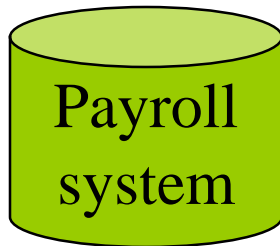
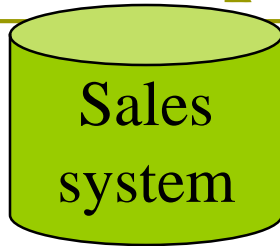
- **Integrarea** asigură conectivitatea între fluxurile de procese economice funcționale



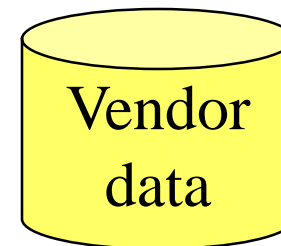
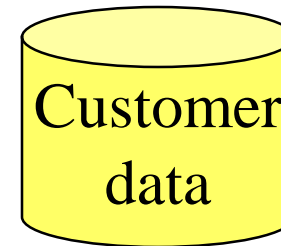
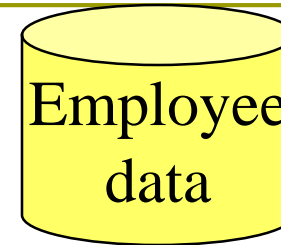
- **Funcționalitatea** a unui sistem ERP asigură fluxurile de procese economice din cadrul fiecărei funcțiuni

# Orientare pe procese/ pe subiecte

---



ERP



DW

# BI si ERP

---

- ERP orientarea pe **processe economice**
- DW orientarea pe **subiecte**
- ERP -BD unica, imensa, cu mii de tabele, care nu se preteaza pentru interogari ad-hoc si analize complexe
- ERP – **avantaj** pentru proiectarea și implementarea DW
  - **omogenitatea** sistemelor sursă și, implicit,
  - modalități mult mai facile de **achiziție** a datelor și de **asigurare a calității**
  - posibilitatea **consolidării datelor** la nivel de companie în cazul firmelor cu mai multe filiale

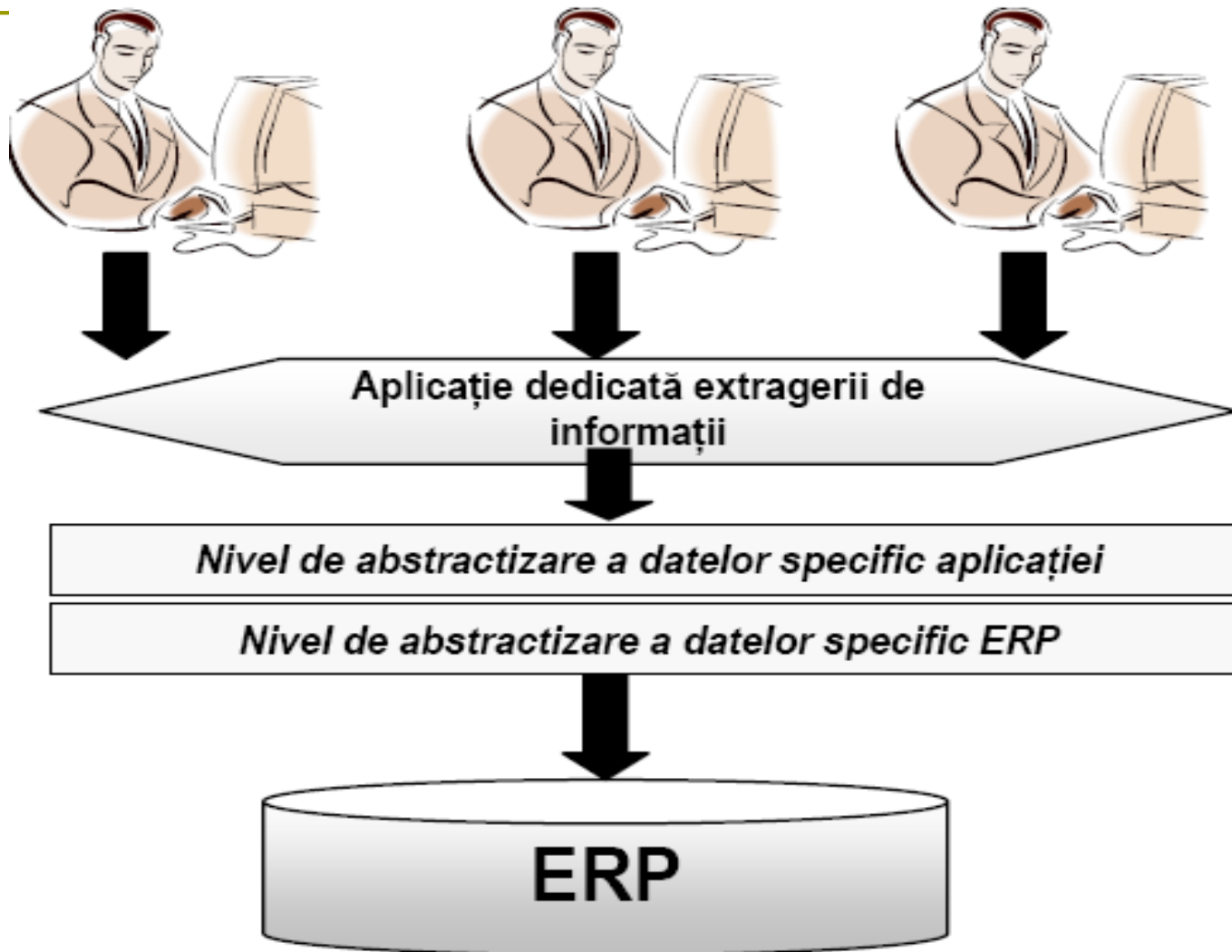


# a. Arhitecturi: Sistem BI cu acces direct la datele din sistemul ERP

---

- ❑ integrat prin intermediul unor aplicații specifice de interogare a datelor.
- ❑ suprapun peste primul nivel de abstractizare al modelului ERP **un nivel de abstractizare propriu**, specific fiecărui utilizator
- ❑ sunt realizate **interfețe dedicate** fiecărui modul din sistemul integrat.
- ❑ **Dezavantaje**
  - limitele impuse de suporturile tehnice.
  - viziunea istorică se suprapune rareori cu necesitățile sistemelor tranzacționale
  - este o soluție de compromis ce poată fi exploatată temporar

# Acces direct la datele ERP

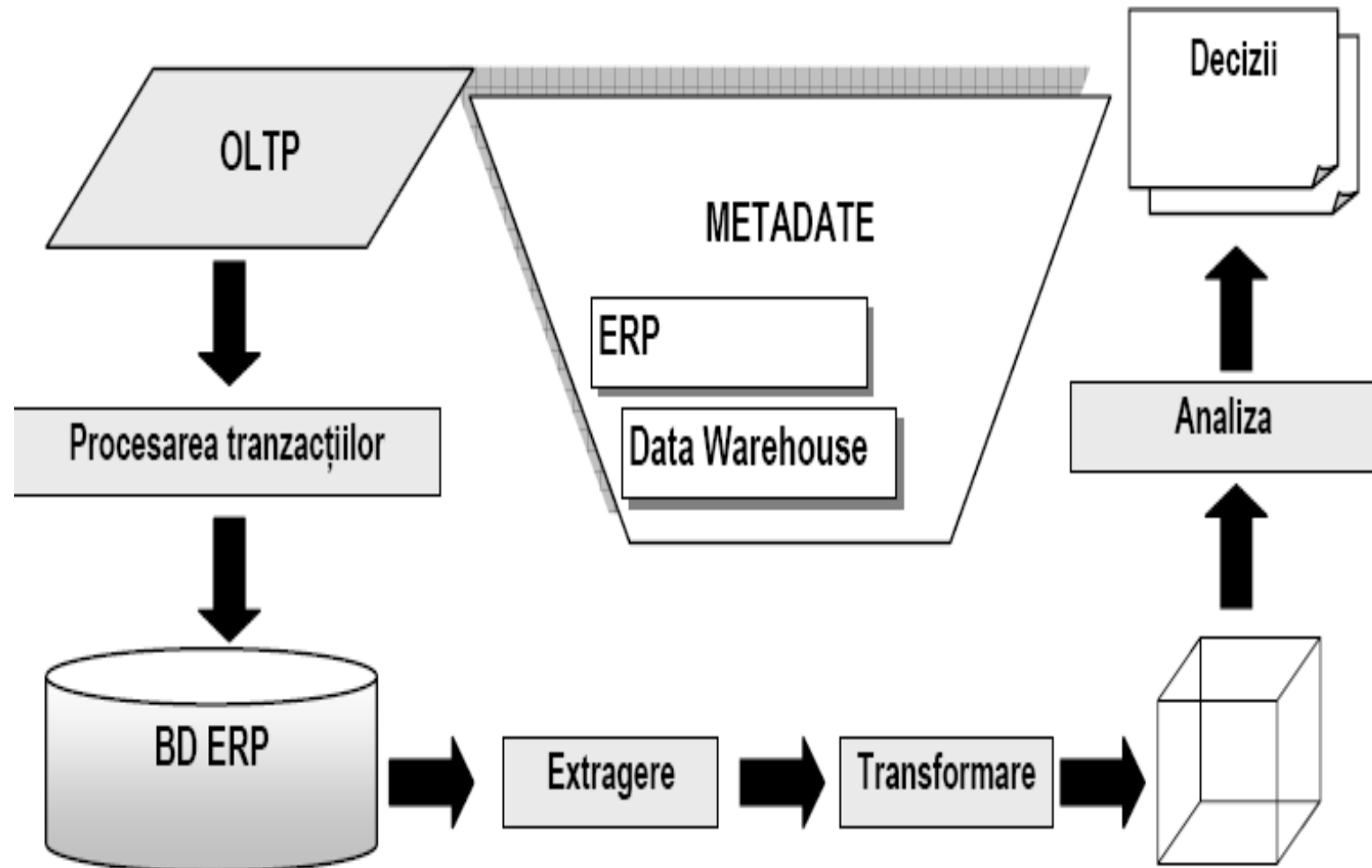


## b.Arhitecturi: Depozit de date atașat ERP

---

- ❑ sistem de asistare a deciziei **specializat**, construit pe baza unui **depozit de date** sau a unei colecții de **data marts**.
- ❑ **dicționar de date propriu**
- ❑ ca aplicație independentă sau ca un modul al ERP (SAP BI)
- ❑ eforturi considerabil mai mari atât în etapele de proiectare și implementare => avantaje prin prisma performanțelor în exploatare

# Depozit de date atașat ERP



# SAP BI (Business Information Warehouse)

---

- **Business Content** = container ce cuprinde
  - Infocuburi (peste 420),
  - Query-uri (peste 1700),
  - Rapoarte si
  - Roluri utilizatorcu specific industrial si functional= solutii preconfigurate pentru diferite industrii
- **Extractori** („plug-in“) - extragerea datelor din SAP ERP si incarcarea in SAP BW
  - complet (full extraction) sau
  - partial (delta extraction).