



# Tehnologii informatice de integrare a datelor

Curs 5

Martie 2018



# Tehnologii de integrare

- 1. Depozite de date**
2. Migrarea datelor
3. Baze de date distribuite



# Business Intelligence

- Principala valoare = capacitatea de **a transforma datele în informații.**
- Cine nu obtine informații suficient de repede și nu le prelucrează eficient, rămâne în urmă sau dispare, într-un mediu concurențial din ce în ce mai agresiv.
- **Business Intelligence** se refera la sisteme informatice de identificare, extragere si analizare a datelor disponibile intr-o companie, sisteme al caror **scop** este **de a oferi un suport real pentru luarea deciziilor de business.**
- O soluție de Business Intelligence integrează **datele curente** ale afacerii dar și **date prealabile**, provenind din mai multe programe și aplicații și le consolideaza într-o singură **bază de date optimizată pentru regăsirea și analiza informației.**

# Analiza datelor



- **Cerinte preliminare:** integrarea datelor
- **Analiza datelor:**
  - Inspectarea, curatarea, transformarea datelor pentru a extrage cunostinte utile
  - Transformarea datelor in informatii si oferirea de raspunsuri descriptive unor intrebari predefinite
- **Data mining**
  - Utilizeaza modelarea datelor pentru descoperirea cunostintelor
- **Business intelligence**
  - Se bazeaza pe analiza datelor si data mining
  - Transforma informatiile si cunostintele in actiuni inteligente
  - Se bazeaza pe diferite instrumente de analiza si pe inteligenta artificiala
- **Instrumente:**
  - Statistica, instrumente de modelare a datelor si de simulare

# Niveluri ale BI



- **BI 1.0**

- Implica mai ales manipularea si prezentarea datelor
- Oferă o platforma care permite examinarea datelor pentru a oferi informatii necesare pentru luarea deciziilor

- **BI 2.0**

- Acces in timp real
- Analize profunde cu instrumente avansate cum sunt score card-uri, metrici KPI, cuburi, panouri de bord
- Personalizare

- **BI 3.0**

- Relevanta – adauga la big data informatii de context pentru a imbunatati relevanta analizelor
- Big data – trateaza volume mari de date nestructurate
- Flexibilitate ridicata – nu se mai folosesc cuburile preconstruite, ci se decide dinamic care sunt elementele necesare: cube-less BI
- Sintetice, creative, raspund intrebarilor deschise

# Depozite de date



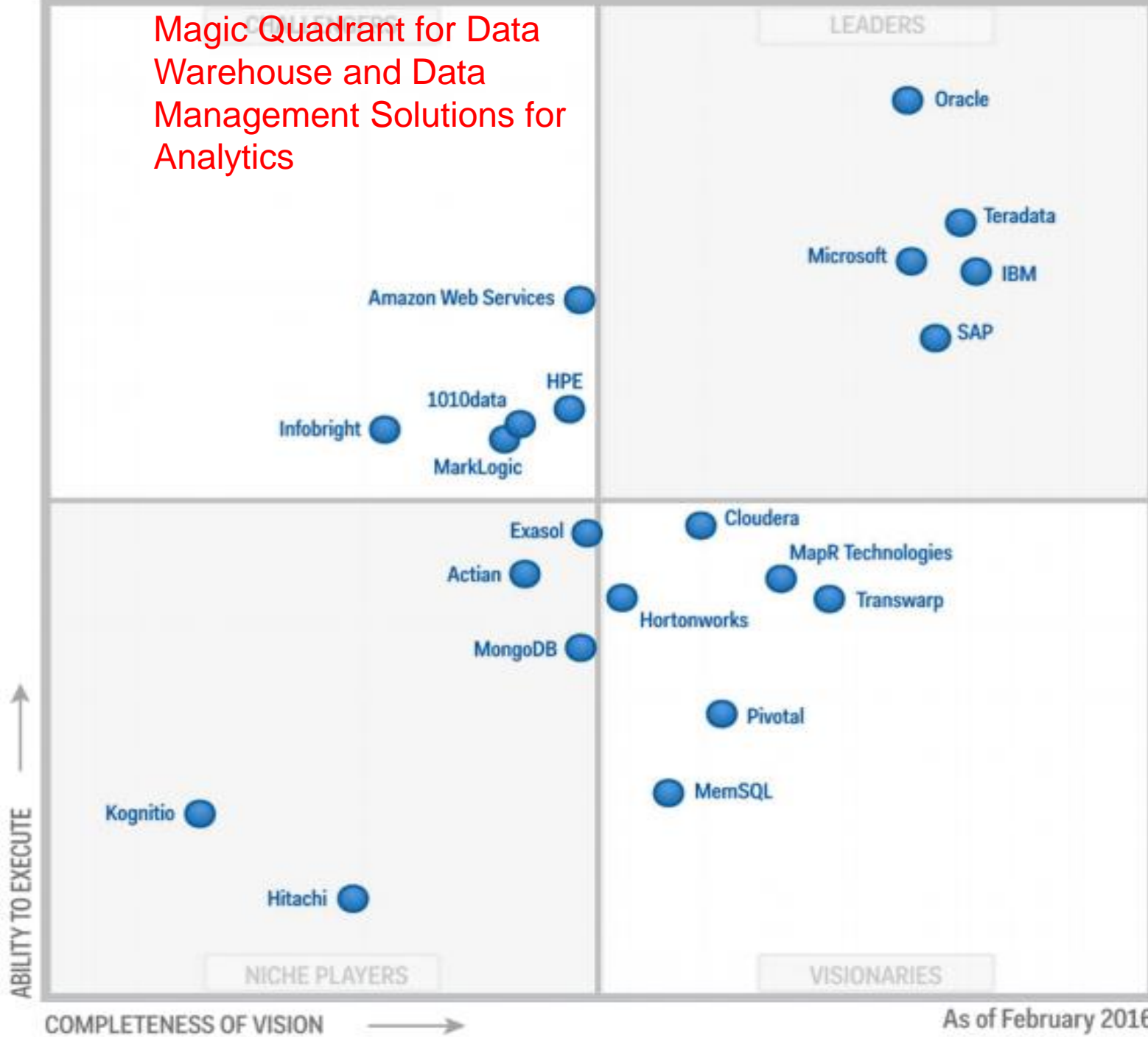
- **Consiliul OLAP 1995:** o stocare centralizată a datelor detaliate provenite din toate sursele relevante din cadrul unei organizații ce permite interogarea dinamică și analiza detaliată a tuturor informațiilor.
- **William Inmon:** o colecție de date orientate pe subiecte, integrate, istorice și nevolatile destinată sprijinirii procesului de luare a deciziilor manageriale
- **+instrumente de interogare, analiza și prezentare a informațiilor**
  - Instrum. de analiza on-line- OLAP
  - Instrum. de data mining



In cifre...

- **Dimensiunea DW – TBytes**
- **Costul implementarii – peste 1 mil \$**
  - Servicii profesionale
  - Software pentru extragere, transformarea, incarcarea si analiza datelor
  - Sisteme hardware si stocarea datelor
- **Gartner Group estimează o creștere dublă pe piața depozitelor de date în raport cu creșterea globală a pieței de IT**

# Magic Quadrant for Data Warehouse and Data Management Solutions for Analytics



As of February 2016



# Scopul DW – integrarea datelor

- **Integrarea datelor**

- modalități unice de codificare, sistem de unități de măsură consistente,
- sistem stabil de reprezentare fizică a datelor,
- convenții clare privind modul de reprezentare a datelor calendaristice,
- convenții unice privind denumirile datelor.

- **FLEXIBILITATE** – sa se conecteze la niv. intregii organizatii a.i. servere de la furnizori diferiti sa se poata conecta la depozitul existent

- **ARHITECTURA** – adaptare usoara la modificarile de performante, capacitate si conectivitate

- **Data mart** - >500 GB, <1mil \$, <3 luni

# Ce este un depozit de date?

- **William Inmon:** este o colecție de date orientate pe subiecte, integrate, istorice și nevolatile, fiind destinat fundamentării deciziei manageriale.
- O BD pentru luarea deciziilor, **separata** de BD operationala a companiei
- Oferă suport pentru **procesarea informațiilor**, oferind o platforma de **date istorice consolidate** pentru analiza
- Structurile de date într-un depozit de date sunt optimizate pentru o **regasire** și o **analiza rapida**.

## a.DW – orientat pe subiecte (Inmon)

- Organizat pe subiecte importante: **client, produs, vanzari.**
- Accent pe modelarea si analiza datelor de catre decidenti
- Oferă o **perspectiva simpla si concisa** asupra anumitor subiecte, **excluzand datele care nu sunt utile in procesul de luare a deciziilor**



## b.DW – integrat (Inmon)

- Integreaza surse de date multiple
- Tehnici de curatare si integrare a datelor.
- Consistenta in
  - conventiile de numire,
  - structura codurilor,
  - unitatile de masura folosite de diferitele surse
    - E.g., Pret hotel: moneda, taxe, mic dejun inclus, etc.

## c.DW – istorice (Inmon)

- **Datele sunt istorice** și sunt actualizate la intervale regulate.
- Orizontul de timp este mult mai mare decât la sist. operationale (ex: 5-10 ani)
- Fiecare element structural cheie al depozitului:
  - Contine o **referire temporală**, implicită sau explicită, ceea ce nu are loc la datele operationale

## d.DW - nevolatii (Inmon)

- Un depozit separat fizic de date transformate din mediul operational
- In DW **nu au loc actualizari operationale ale datelor.**
  - Nu necesita mecanisme de procesarea tranzactiilor, recuperare si controlul concurentei
  - Sunt necesare doar 2 operatii pentru accesarea datelor :
    - **Incarcarea datelor si accesul la date.**
- **Actualizare** doar adăugarea periodică a unor date extrase din sistemele operationale
- Preocupare pt. **optimizarea accesului la date**: denormalizare, sumarizare, statistici ale accesării și reorganizare dinamică a indexării

# Aplicatii ale depozitelor de date

- **Telecomunicatiile.**

- folosirea rețelei,
- profilul clientilor care folosesc un anumit serviciu,
- profitabilitatea produselor si serviciilor oferite.

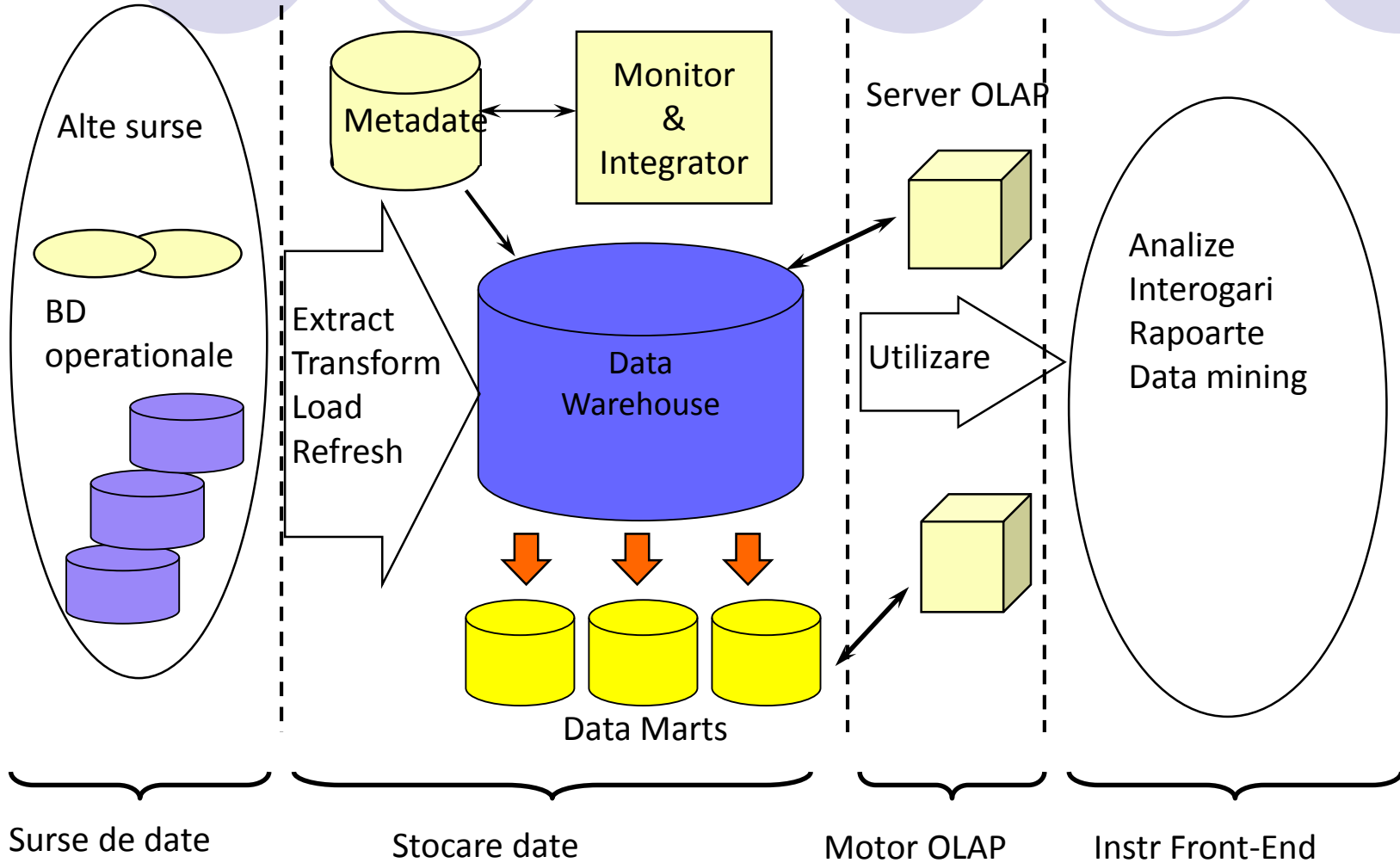
- **Bancile**

- gestionarea profitabilitatii generale, prin analizarea profitabilitatii pe produs si pe client.
- determinarea profilurilor clientilor pentru a directiona cât mai eficient campaniile de marketing.

- **Comertul cu amanuntul.**

- trendul vânzarilor în functie de anotimp, vacante, campanii de publicitate, activitatea competitorilor.
- mentalitatile si obiceiurile cumparatorilor = intrari in sistemul de dirijare a actiunilor promotionale si a altor campanii de marketing
- analiza trendului performantelor
- vânzari încrucisate
- profilul consumatorului si piata tinta.

# Architectura multinivel





# De ce un depozit de date separat?

- **Performante mai bune**

- **SGBD**— potrivit pt OLTP: metode de acces, indexari, controlul concurenței, recuperare.
- **Depozit** —potrivit pt OLAP: cereri complexe, perspective multidimensionale, consolidare

- **Functii si date diferite**

- Date: luarea deciziilor necesita date istorice
- Consolidarea datelor: luarea deciziilor necesita consolidari de date din surse eterogene
- Calitatea datelor: datele din surse diferite au reprezentari, codificari si formate diferite care trebuie reconciliate

# Tipuri de DW

## 1 DEPOZITE DE ÎNTREPRINDERE (ENTERPRISE WAREHOUSE)

- întreaga structură organizațională
- un volum extins de date: atât informații **detaliate**, cât și **agregate**.
- suporturi hardware performante.
- costurile și timpul de proiectare și implementare sunt considerabile,

## 2 DATA MART

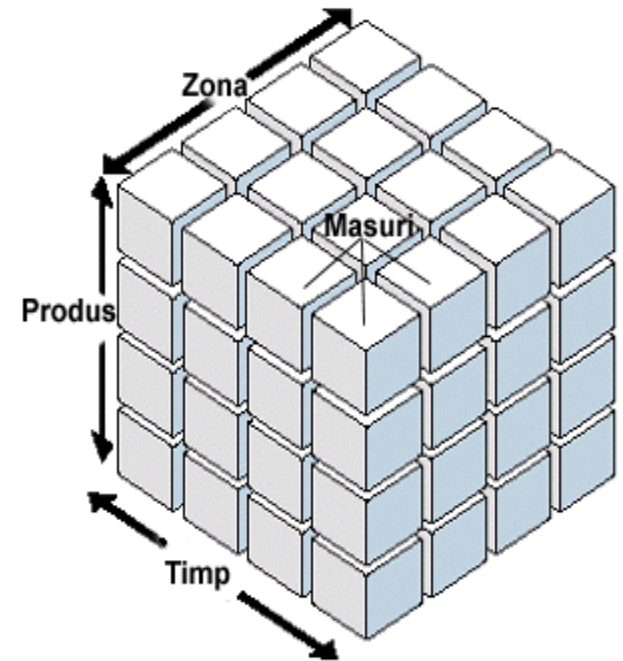
- specific unui anumit subset de cerințe sau unui departament din cadrul organizației
- de regulă, datele conținute într-un data mart sunt **agregate**.
- costurile și termenele de implementare sunt considerabil reduse

## 3 DEPOZIT VIRTUAL (VIRTUAL WAREHOUSE)

- o serie de vederi (*views*) realizate direct asupra BD operaționale.
- procesele de agregare pot afecta capacitățile de prelucrare ale serverelor utilizate în activitatea operațională,
- aparent ușor de implementat, necesită capacități de procesare deosebite.
- necesită curățare și consolidare în timpul rularii

# Modelul multidimensional

- permite vizualizarea datelor prin mai multe filtre sau **dimensiuni** in acelasi timp.
- Dimensiuni=coordonate=  
categorii de informație.
- De ex:
  - Care sunt vanzarile reale in comparatie cu cele previzionate pe zona, pe vanzator, pe produs?
  - Care este profitabilitatea pe produs, pe client?



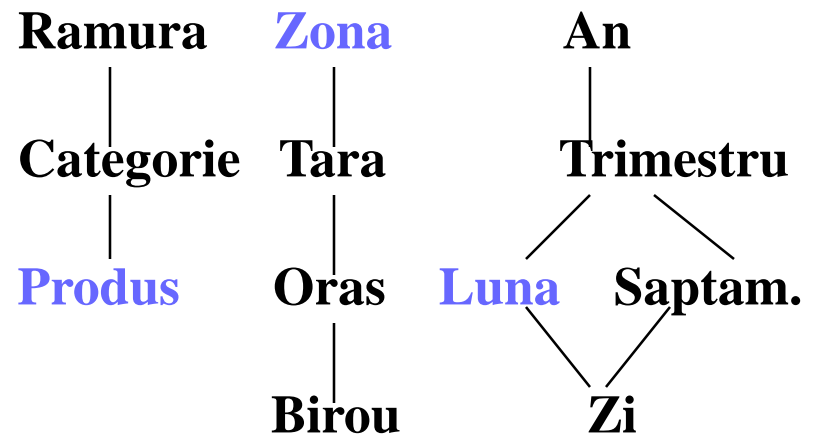
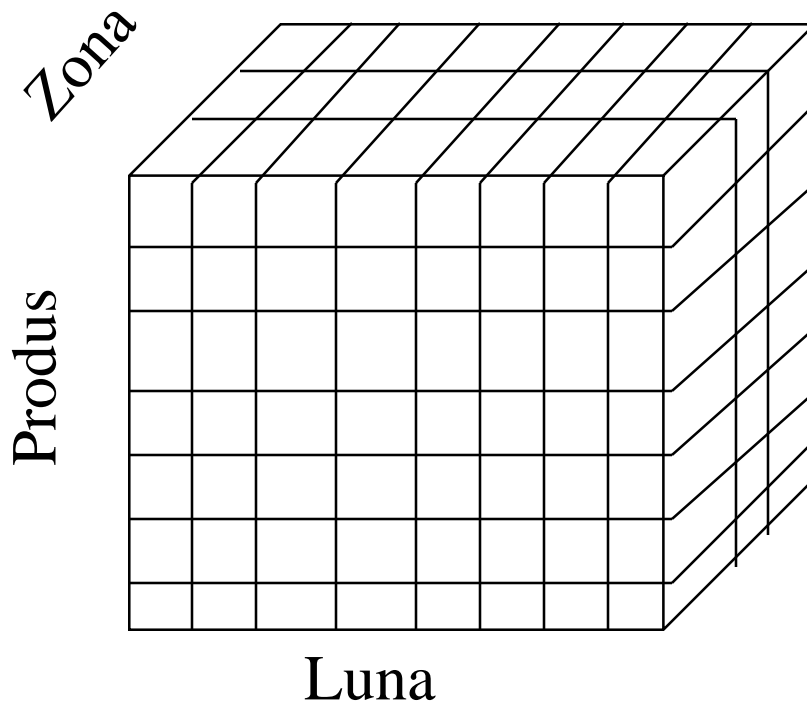
# Obiecte DW

- **Tabelele de fapte** (masuri)
  - conțin faptele și cheile externe către tabelele de dimensiuni.
  - de obicei date numerice - totalizate și analizate pe diferite niveluri.
- **Tabele dimensiuni**
  - categorii de informații care organizează datele
  - fiecare tabelă dimensiune are câte o cheie principală
  - câmpurile sunt de obicei textuale - **sursă pentru restricții și pentru rândurile din rapoarte.**
  - datele sunt de obicei colectate la nivelul cel mai de jos și mai detaliat și agregate pe nivelele superioare pentru analiză.
- **Atribut** - un nivel al unei dimensiuni, într-o **IERARHIE**
- **Ierarhiile**
  - sunt structuri logice utilizate pentru ordonarea nivelelor de reprezentare a datelor.
  - definesc caile de navigare în interiorul datelor, permițând detalierea graduală a datelor.

# Date multidimensionale

- Volumul vanzarilor – functie de produs, luna, si zona

Dimensiuni: Produs, Zona, Timp  
Ierarhii:



# Exemplu: Vanzari de fructe

<b>Timp</b>	<b>Suma</b>
Trim 1	16000
Trim 2	16000
<b>Total Timp</b>	<b>32000</b>

<b>Piata</b>	<b>Suma</b>
Brasov	8000
Sibiu	8000
Arad	8000
Iasi	8000
<b>Total Piata</b>	<b>32000</b>

<b>Produs</b>	<b>Suma</b>
Mere	8000
Cirese	8000
Struguri	8000
Pepeni	8000
<b>Total Produs</b>	<b>32000</b>

		<b>Brasov</b>	<b>Sibiu</b>	<b>Arad</b>	<b>Iasi</b>	<b>Total</b>
<b>Trim. 1</b>	Mere	-	-	2500	1500	<b>4000</b>
	Cirese	-	-	2000	2000	<b>4000</b>
	Struguri	1000	3000	-	-	<b>4000</b>
	Pepeni	2000	2000	-	-	<b>4000</b>
	<b>Total trim 1</b>	<b>3000</b>	<b>5000</b>	<b>4500</b>	<b>3500</b>	<b>16000</b>
<b>Trim 2</b>	Mere	4000	-	-	-	<b>4000</b>
	Cirese	1000	3000	-	-	<b>4000</b>
	Struguri	-	-	1500	2500	<b>4000</b>
	Pepeni	-	-	2000	2000	<b>4000</b>
	<b>Total trim 2</b>	<b>5000</b>	<b>3000</b>	<b>3500</b>	<b>4500</b>	<b>16000</b>
	<b>Total</b>	<b>8000</b>	<b>8000</b>	<b>8000</b>	<b>8000</b>	<b>32000</b>

# Agregari si granularitate

- **Granularitatea** – reprezinta nivelul de detaliere la care sunt pastrate datele in depozit
- In functie de cerintele de analiza, datele se pot pastra la nivel mai detaliat sau mai agregat (depinde de niv. de detaliere a dimensiunilor)
- **Agregarea** datelor- cresterea performantelor DW
- 10 magazine, 100 produse/marca, vanzari saptamanale

<b>Dacă o interogare necesită...</b>	<b>Atunci trebuie parcurse</b>
1 Produs, 1 Magazin, 1 Săptămână	doar 1 înregistrare din schemă
1 Produs, Toate magazinele, 1 Săptămână	10 înregistrări din schemă
1 Marcă, 1 Magazin, 1 Săptămână	100 înregistrări din schemă
1 Marcă, Toate magazinele, 1 An	52.000 înregistrări din schemă

# De la relational la multidimensional

- premise diferite, tehnici diferite și produc BD cu structuri diferite.
- modul de abordare a datelor (utilizator/date):
  - **model multidimensional** - dimensiuni cât mai apropiate de cele naturale și de **perspectiva utilizatorului**.
  - **model relational** – perspectiva datelor
- model multidimensional:
  - o BD mult **mai ușor de consultat și de interogată** la un nivel înalt, sintetic, agregat
  - o BD cu **mai puține tabele și chei** de administrat decât modelul relational



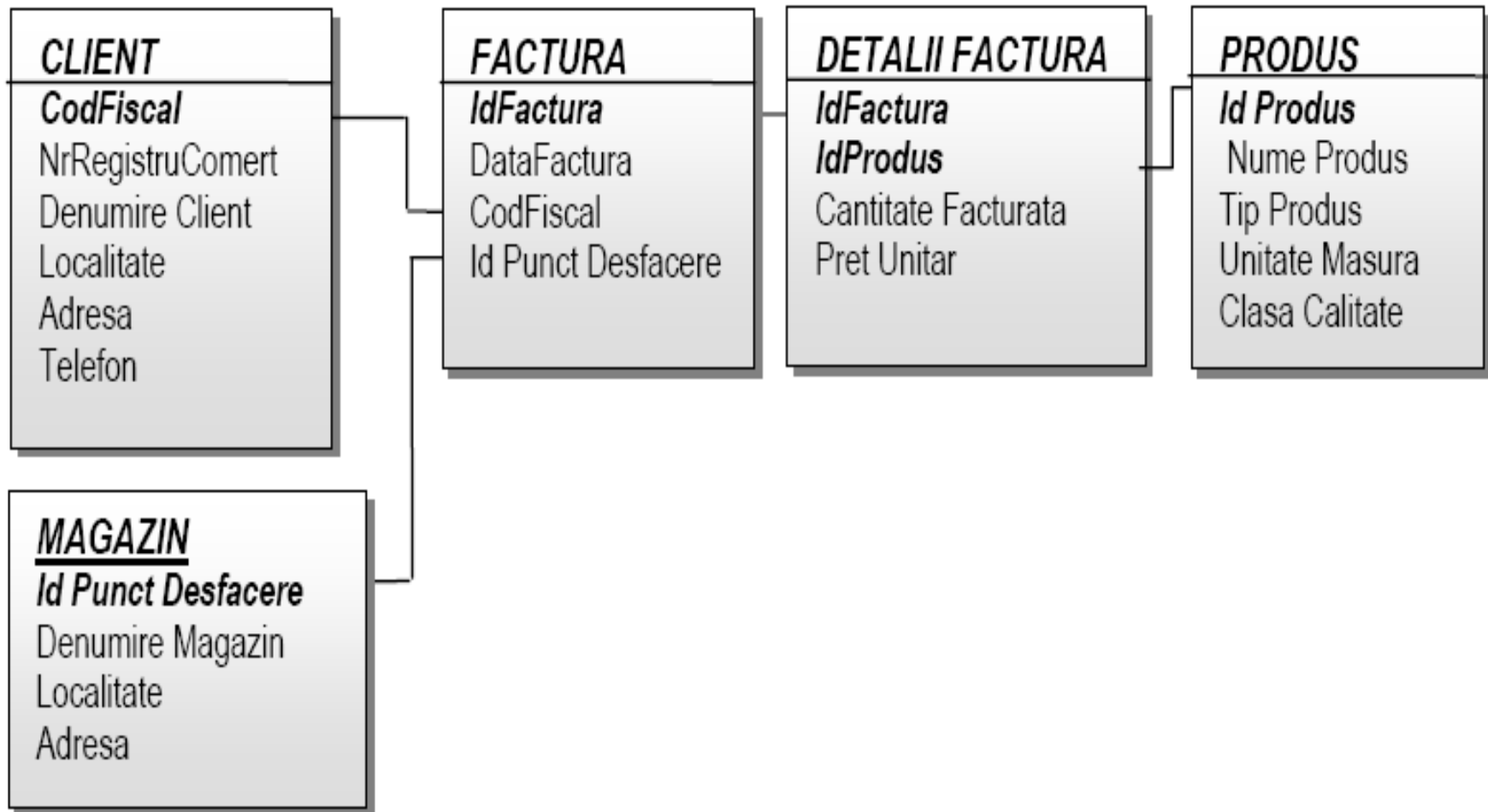
# Normalizare in BD, redundante in DW

- procesul de transformare succesivă a unei BDR în vederea aducerii sale într-o formă standard optimizată
- eliminarea anomaliilor, redundanțelor, dependențelor nedorite între date
- **Anomalii de actualizare**
  - limitarea posibilităților de inserare a datelor
  - pierderi de date la ștergere
  - apariția de inconsistențe la modificarea datelor
- **Dependente**
  - **Dependență funcțională** – A depinde funcțional de un B dintr-o tabelă dacă fiecărei valori a lui A îi corespunde numai o valoare a lui B. B **depinde funcțional complet** de un grup de attribute dacă B este dependent funcțional de fiecare atribut din grup.
  - **Dependentă tranzitivă** –daca B depinde de A și C depinde de B atunci C se află în dependență tranzitivă față de A.
  - **Dependență multivaloare** – dacă valorii unui atribut A îi corespund două sau mai multe valori ale atributului B

# Formele normale

- **Forma normală 1 (FN1)** *dacă* attributele sunt la nivel **atomic** și au fost eliminate **grupurile de attribute** repetitive
- **Forma normală 2 (FN2)** *dacă* este în FN1 și nu există **dependențe funcționale parțiale** pentru attributele non-cheie
- **Forma normală 3 (FN3)** *dacă* este în FN2 și nu există **dependențe funcționale tranzitive** pentru attributele non-cheie
- **Forma normală 4 (FN4)** *dacă* este în FN3 și există cel mult o dependență funcțională multivaloare pentru attributele non-cheie
- **Forma normală 5 (FN5)** *dacă* este în FN4 și nu există dependențe joncțiune pentru attributele non-cheie

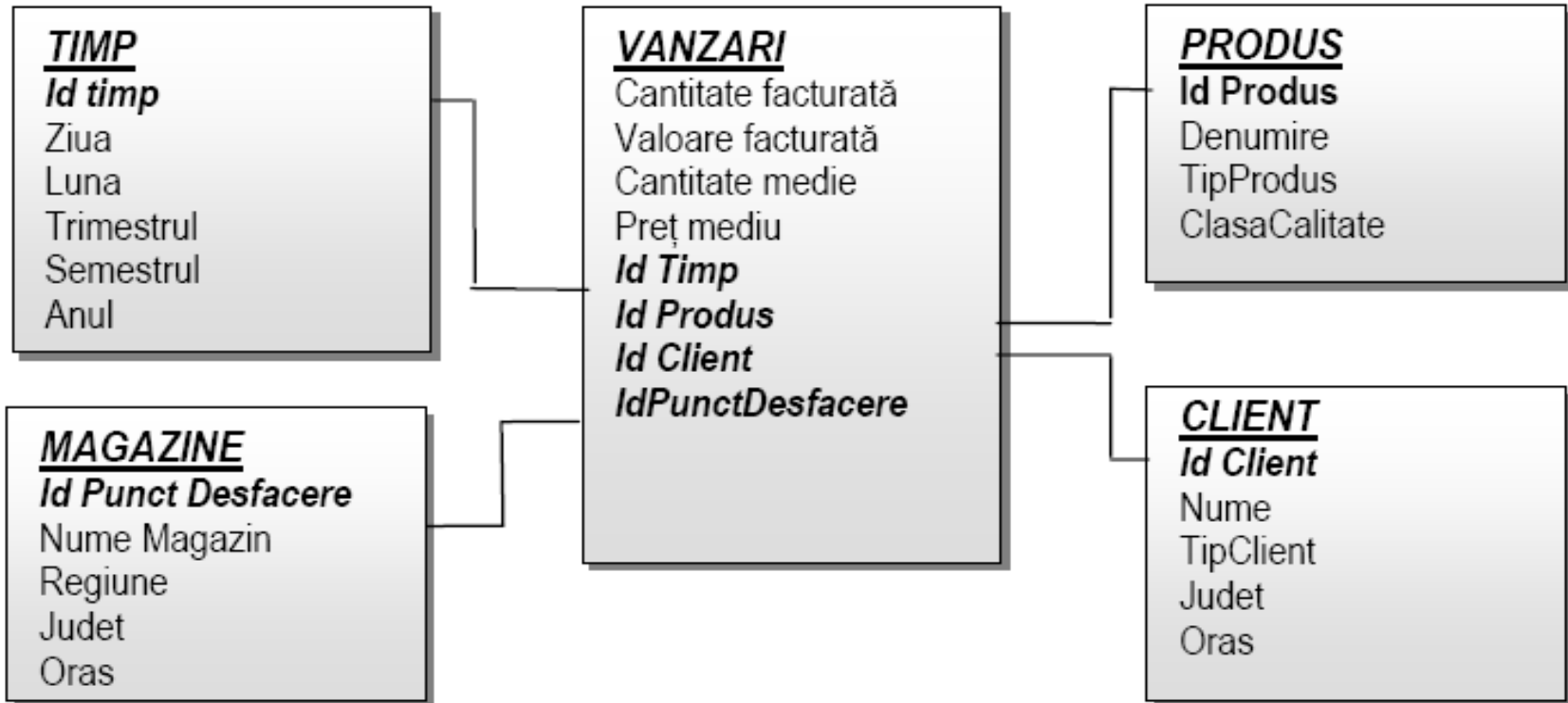
# Evidenta facturi – model relational



# a. Structura DW – Schema STEA

- cel mai des utilizat model de organizare al depozitelor de date
- **tabela de fapte** cuprinde, fără redundanțe, marea parte a datelor
- tabela de fapte este conectata la tabelele dimensiune pe baza cheilor externe pe care acestea le conțin.
- **star join** = legatura stabilita între un tabel de fapte si tabelele dimensiune
- **star query** = jonctiunea dintre un tabel de fapte si mai multe tabele dimensiune
- **Avantaj**: performante optime pentru interogariile dintr-un depozit de date

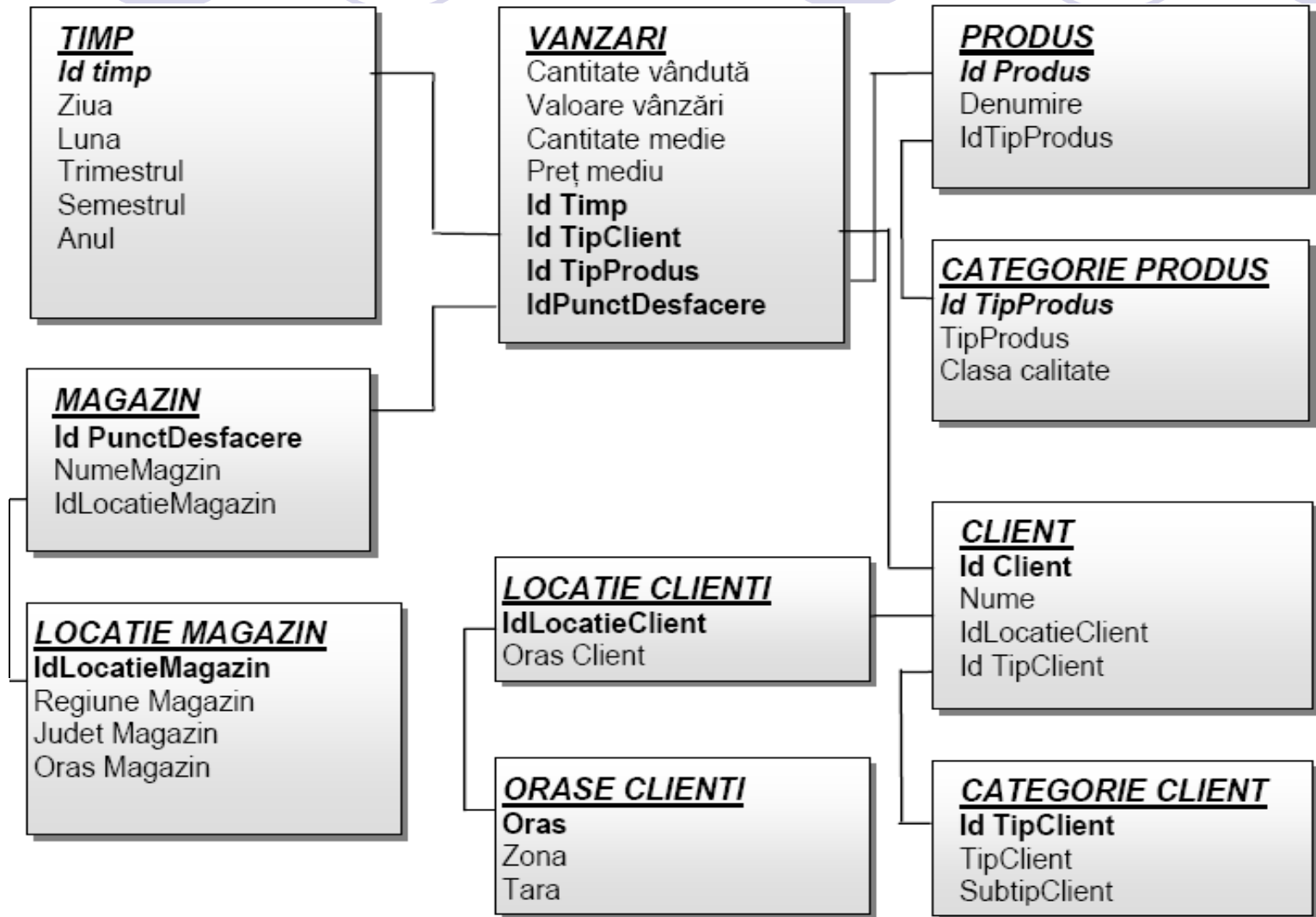
# Ex: Schema STEA



## b. Structura DW– Schema fulg de zapada

- “**seminormalizat**”, avantajele modelului relațional.
- tabelele dimensiune respecta regulile de normalizare din modelul relațional => economie de spațiu
- nu va conduce la reducerea spațiului pt tabela de fapte
- **Avantaje:**
  - Redundanta redusă
  - Usor de întreținut
- **Dezavantaje:** la cereri de interogare complexe(join)=> crește timpul de răspuns

# Ex: Schema fulg de zapada

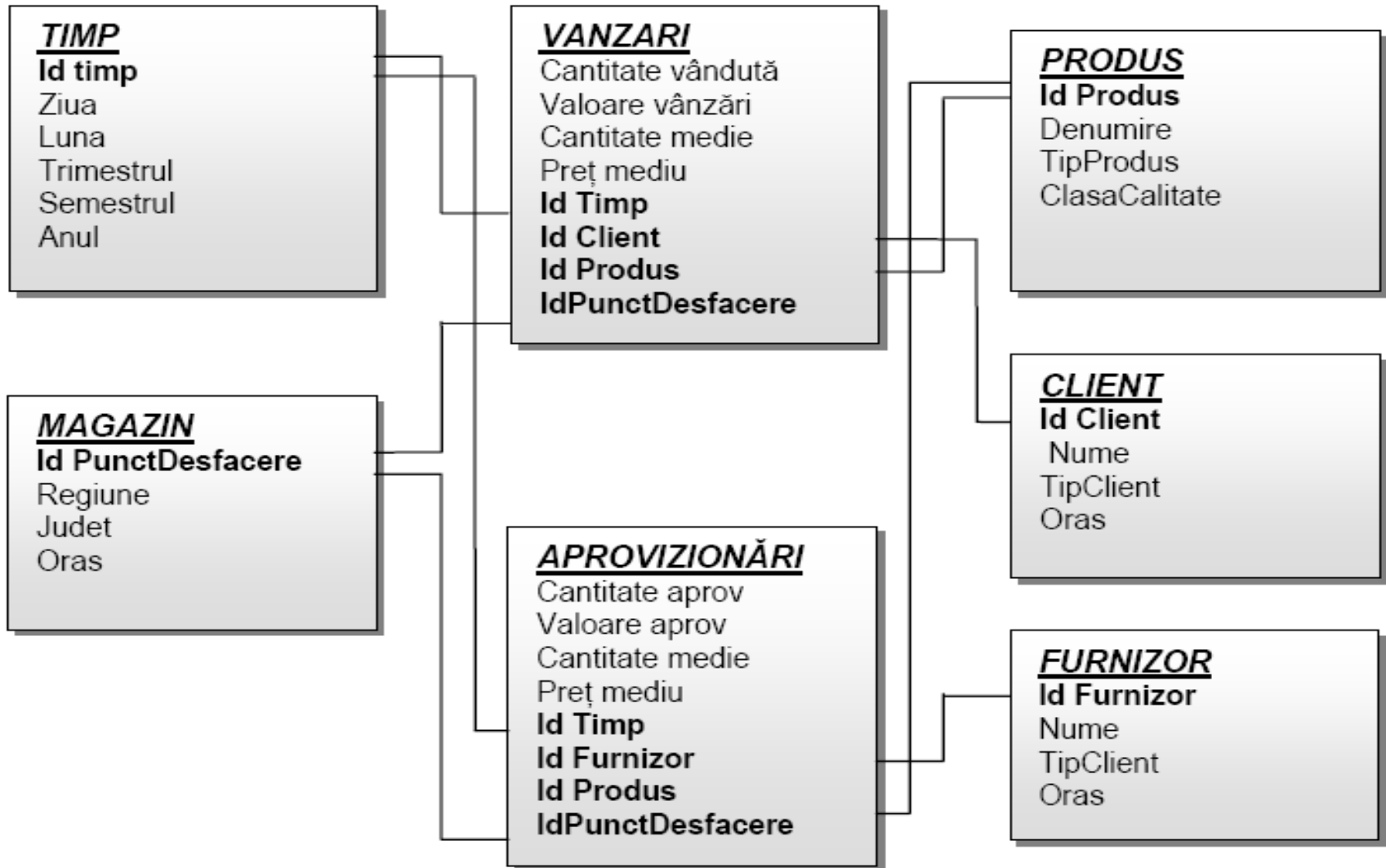


## c. Structura DW – Schema constelație de fapte

- Schema galaxie
- mai multe tabele de fapte, conectate ce utilizează aceleași tabele-dimensiune
- pe lângă tabela de fapte **Vânzări**, o tabelă suplimentară de fapte **Aprovizionări**, legata de dimensiuni



# Ex: Schema constelație de fapte



# Paralela între prelucrarea relatională și cea analitică

Caracteristici	Modelul relational	Modelul multidimensional
Organizarea datelor	Tabela	Dimensiuni, tabele de fapte, cub de date
Nivelul datelor	Detaliu	Agregat
Operația tipică	Actualizare	Raportare și analiză
Nivelul de analiză cerut	Scazut	Ridicat
Volum de date per tranzacție	Redus	Mare
Vârsta datelor	Curente	Istorice, curente, previzionate



# Piata job-urilor

- Se previzioneaza o lipsa mare de personal in urmatoorii 5-10 ani pe piata analizei de date
- O piata foarte dinmica
- Acumulari tot mai mari de date
- Noi tehnologii si instrumente
- E nevoie atat de instrumentele traditionale de analytics dar si de expertiza tehnica pentru date nestructurate
  - *IBM are aproape 10000 de consultanti analisti de business si 400 matematicieni*

# Nevoia de analiza a datelor

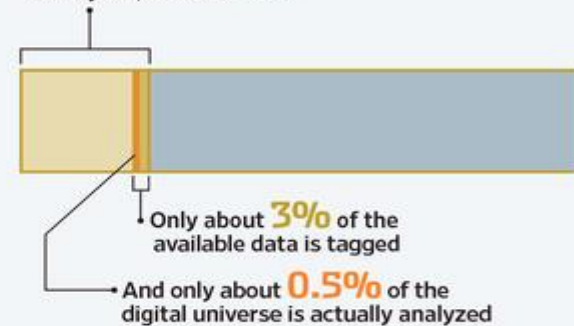
## Bigger and Bigger Vast Potential

■ From now until 2020, the digital universe—the total amount of data created, replicated and consumed each year all over the world—is expected by International Data Corp. to grow steadily to more than 14 times its current size

■ Data in millions of petabytes



23% of this digital universe is information that might be valuable if analyzed, IDC estimates



## Missing Out

■ EMC Corp. describes business-intelligence analysts as those focused on using corporate data to monitor or manage various business concerns, and data scientists as those who use advanced analytical tools to tease predictive insights and product innovations out of data.

**38%** of business-intelligence analysts and data scientists strongly agree that their company uses data to learn more about customers.

**22%** of data scientists strongly believe their company gives all employees access to the data needed to run experiments that could result in important innovations.

The most commonly cited barriers to data-science adoption include:

**32%**  
Lack of skills or training

**32%**  
Budget constraints

**14%**  
Organizational structure

## Help on the Way?

■ Many in the profession expect a persistent shortage of analytical talent.

**83%** of data scientists and business-intelligence analysts think new tools and technology will increase demand for data scientists

**64%** believe demand for data scientists will outpace supply

**58%** believe the best source of new data scientists will be university students; 12% think the best source will be today's business-intelligence analysts

# Tehnologii de integrare



1. Baze de date distribuite
2. Depozite de date
3. **Migrarea datelor**

# 3. Migrarea datelor



- Migrare sau reproiectare la schimbarea BD
- Avantaje reproiectare
  - posibilitatea de a începe de la zero și a elimina slăbiciunile structurale;
  - adoptarea de noi tehnologii;
  - crearea unei fundații proaspete pentru noul sistem
- Dezavantaje reproiectare
  - analiza, proiectarea și implementarea unui nou sistem solicită mult timp și resurse
  - este posibil ca noul sistem să fie mai puțin funcțional decât vechiul

# Factori ce influenteaza migrarea

- Diferențele de **sintaxă SQL** între principalele SGBD-uri;
- Integrarea de **restricții de integritate** și **algoritmi** atât în BD sursă, cât și în destinație
- **Asistent de migrare**, care să automatizeze cele mai multe sarcini, iar administratorul BD să facă doar corecții minore și de finețe.
- Interdependența dintre obiectele BD
- Volumul mare de date – durata mare transfer

# Etapele migrării datelor



- A. Export și conversie
- B. Transfer și procesare
- C. Import



# A. Export si conversie

- Se exporta si se convertesc toate/ o parte din obiectele BD
    - Tabele
    - Viziuni
    - Proceduri/ functii/ pachete stocate
    - Declansatori
  - Redenumiri sau schimbari de tipuri
- => **Fisiere ASCII cu comenzi SQL** pentru crearea structurii si cu date pentru popularea BD

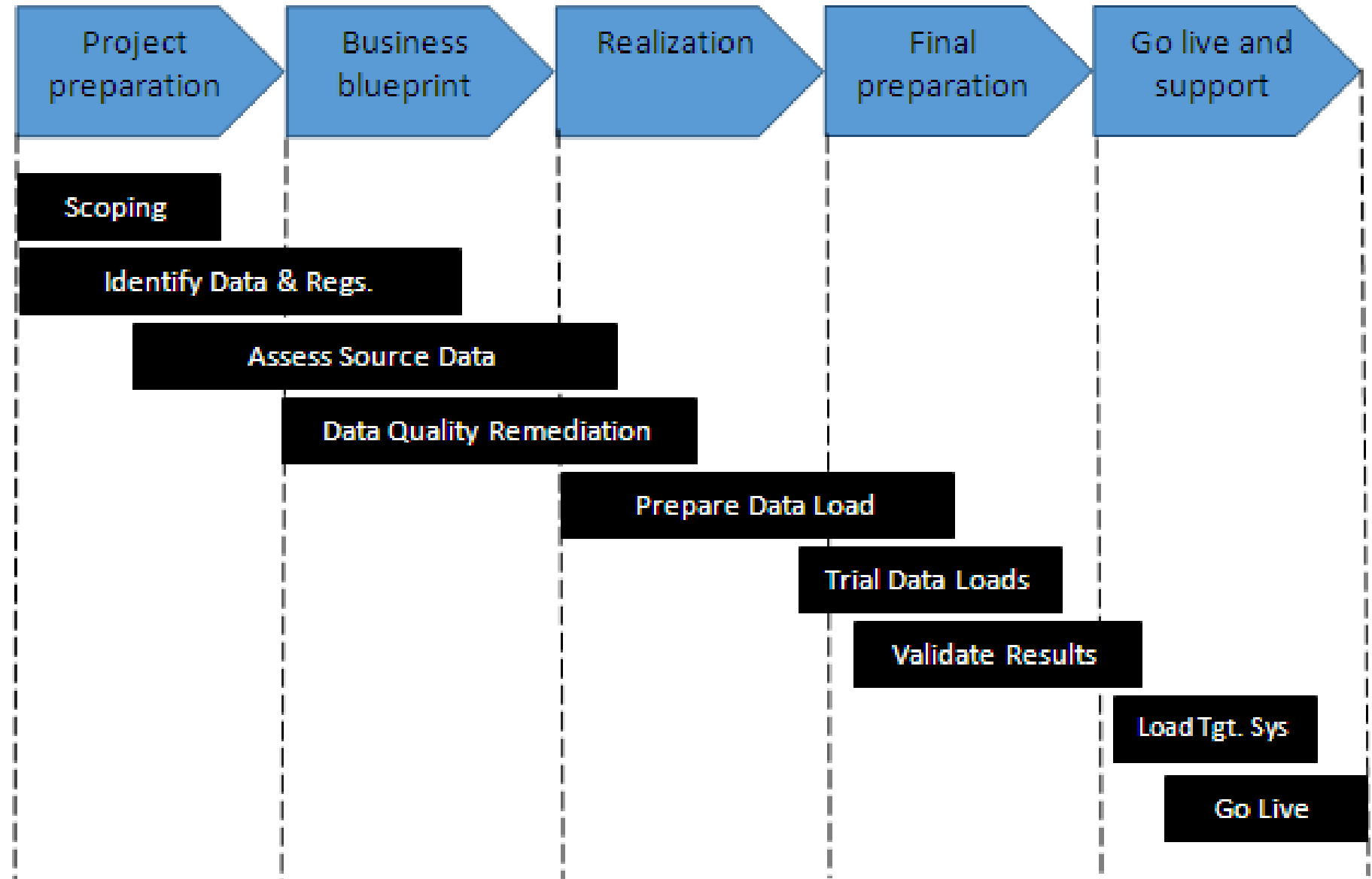
## B. Transfer si procesare scripturi

- Optionala, daca e nevoie de transfer
- Procesarea scripturilor transferate – modificari pt nevoi **neacoperite** de agentul de migrare folosit

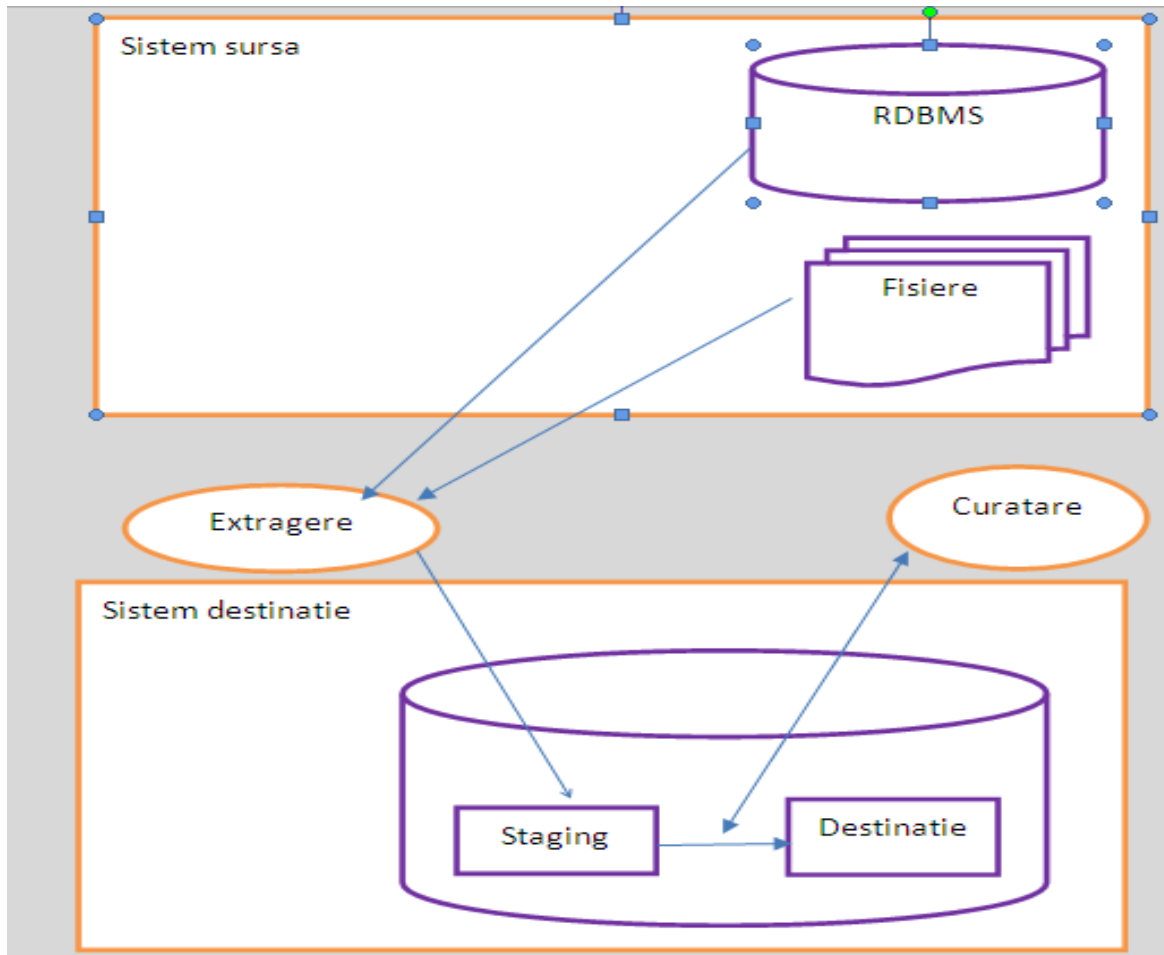
# C. Import

- Scriptul creat la A si prelucrat la B e executat pe BD destinatie
- Utilitare pt executia scripturilor:
  - **SQL Plus** pentru Oracle;
  - **CLP** (Command Line Processor) pentru IBM DB2;
  - **ISQL** pentru Ms SQL Server și SyBase;
  - **linia de comandă MySQL**.
- Utilitare pt. incarcare date din fisiere ASCII:
  - **SQL Loader** pentru Oracle;
  - **LOAD/IMPORT** pentru IBM DB2;
  - **BCP** pentru SQL Server și Sybase;
  - **LOAD DATA INFILE** pentru MySQL;
  - **BUTIL** pentru Persasive SQL.

# Planificarea unui proiect de migrare a datelor

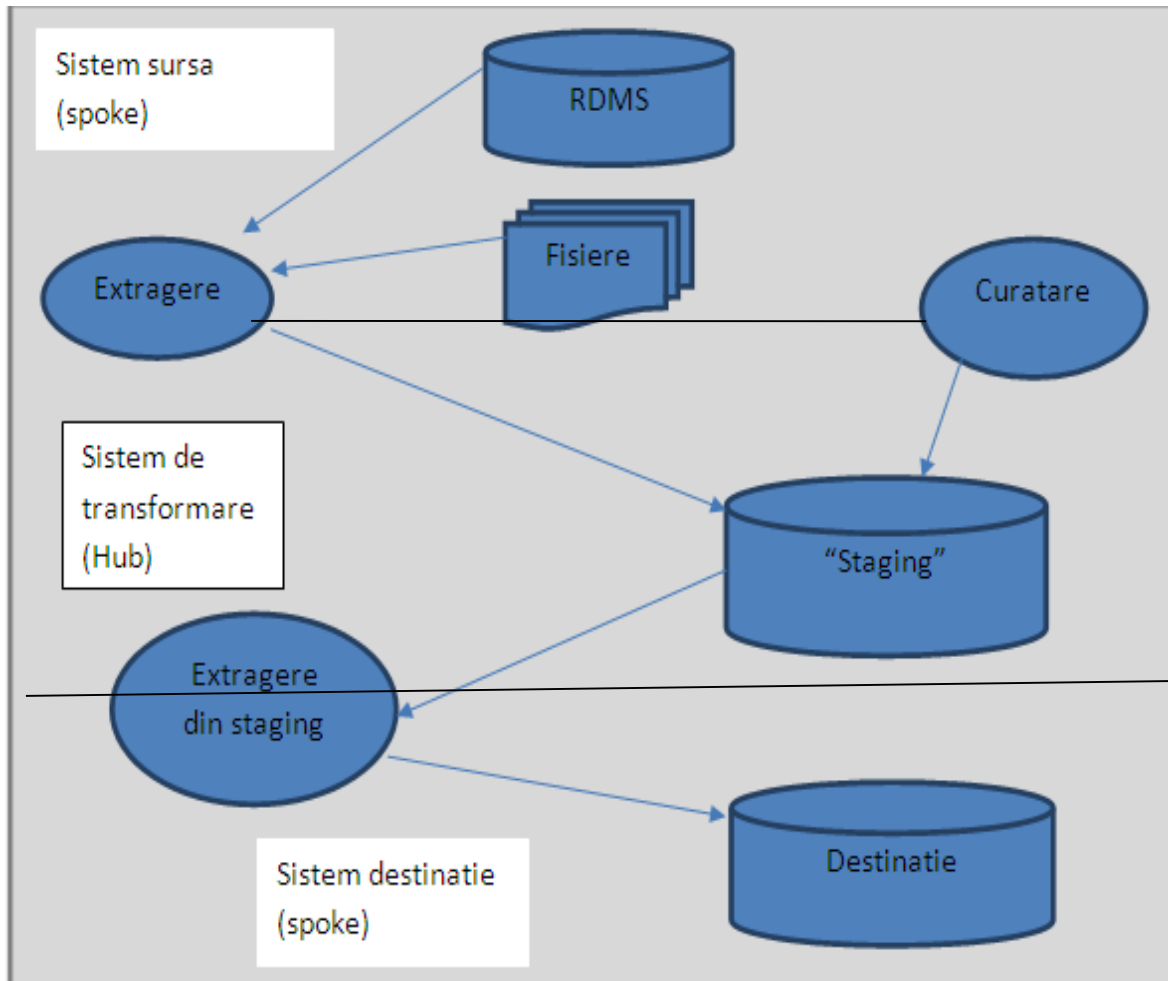


# Arhitectura de migrare punct-la-punct



- Transformările locale înseamnă ca procesul de migrare a datelor s-a terminat, datele au ajuns la sistemul destinatie
- Reduce timpul petrecut “in retea”

# Arhitectura Hub-Spoke de migrare a datelor



- Poate integra orice numar de sisteme sursa/destinatie

- Regulele datelor sunt pastrate la un nivel separat

# Strategii de migrare



- **Strategia “Big Bang”**

- migrarea tuturor datelor si trecerea la noul sistem in acelasi pas.
- avantajul -lipsa nevoii interoperabilitatii intre vechiul sistem si noul sistem.
- dezavantaj- durata mare de “downtime” sau neputinta testarii in productie

- **Strategia “Chicken Little”**

- sistemul sursa este divizat in unitati cu cat mai putine interdependente
- vechiul sistem si noul sistem ruleaza in paralel in timp ce modulele sunt transferate
- migreaza datele incremental, asigurand integritatea informatiei.
- Avantaj: timpul de stabilizare permis intre migrarile modulelor., testare

- **Strategia “Butterfly”**

- sursa este migrata iterativ pana cand diferenta dintre cele doua sisteme a atins pragul prestabilit, moment in care, restul informatiei este transferat si noul sistem este pornit
- nu se foloseste de portalul intre cele doua sisteme.
- sistemul tinta nu se afla in productie deci nu trebuie sincronizat dupa fiecare migrare.