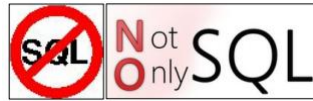




Cap.11. Baze de date Non-SQL

CUPRINS



1. Definiere BD NoSQL
2. Caracteristici
3. Tipuri de BD NoSQL
4. Exemple SGBD NoSQL

1



1. Definiere BD NoSQL

BD NoSQL (Next Generation Databases, Modern web-scale databases) = BD non-relationale, distribuite, open-source, scalare orizontala.

De ce a aparut necesitatea utilizarii unor BD NoSQL?

1. BD relaționale : sunt cel mai mult folosite în aplicații, **DAR**

2. Sunt probleme cu utilizarea lor, ex.:

- BD =mulțime de relații (tabele). Fiecare relație = schemă fixă. **Modificare scheme -> modificarea aplicațiilor-> costuri mari.**
- Dimensiunea foarte mare a fisierelor BD** daca trebuie memorate cantități foarte mari de text + imagini+ audio+ filme
- Unele **informații** (temporare sau persistente) sunt **greu de memorat** în modelul relațional: coșuri de cumpărături, obiecte 3D,etc
- Interogare dificila** a datelor care depind de alte date (din alte baze de date)



2. Generalitati BD NoSQL

1998: probleme de gestiune date au apărut in web 2.0, : Google, Amazon, Yahoo, Facebook, Twitter, etc. Rezolvarea problemelor -> soluții NoSQL.

2009: semnificatia actuala **Not only SQL**

Avantaje si caracteristici:

- memorare volume mari de date (10-100K servere, 1K=100.000). Ex. **Google's "big data"**
- nu există o structură fixă a datelor
- legaturi între date (prin referințe la date din alte baze de date)
- date replicate pe mai multe servere (partajare și replicare)
- interogările se fac rapid
- se pot gestiona probleme la operatii de tip Insert & Update asincrone



2. Generalitati BD NoSQL

Dezavantaje:

- nu există standarde (similar standard SQL BD relaționale)
- nu se asigură consistența bazei de date (de către sistemul de gestiune)
- nu există metode performante pentru protecția datelor
- modelele propuse sunt la primele versiuni
- există posibilități limitate de interogare
- aproape toate sistemele apărute sunt open-source
- există relativ puțini dezvoltatori software pentru NoSQL



2. Generalitati BD NoSQL

Solutii crestere viteza de gestiune BD de dimensiuni mari:

- ❑ "scalare verticală": se adaugă resurse hard suplimentare sau resurse hard cu caracteristici calitative superioare (memorie internă, procesor, hard-disk, etc.) - dar cu costuri ridicate
- ❑ "scalare orizontală": distribuirea și replicarea datelor, efectuarea calculelor în mai multe noduri ale rețele



2. Generalitati BD NoSQL

Unit	Symbol	Bytes	Server	Cost
Kilobyte	KB	1024	PowerEdge T110 II (basic) 8 GB, 3.1 Ghz Quad 4T	\$1,350
Megabyte	MB	1048576	PowerEdge T110 II (basic) 32 GB, 3.4 Ghz Quad 8T	\$12,103
Gigabyte	GB	1073741824	PowerEdge C2100 192 GB, 2 x 3 Ghz	\$19,960
Terabyte	TB	1099511627776	IBM System x3850 X5 2048 GB, 8 x 2.4 Ghz	\$646,605
Petabyte	PB	1125899906842624	Blue Gene/P 14 teraflops, 4096 CPUs	\$1,300,000
Exabyte	EB	1152921504606846976	K Computer (fastest super computer) 10 petaflops, 705,024 cores, 1,377 TB	\$10,000,000 annual operating cost
Zettabyte	ZB	1180591620717411303424		
Yottabyte	YB	1208925819614629174706176		



3. Infrastructura BD NoSQL

Infrastructura BD NoSQL: mai multe servere, conectate într-o rețea. Toate nodurile fizice au același sistem de operare. Într-un nod fizic se pot organiza mai multe noduri virtuale

SBD relaționale: folosesc un sistem tranzacțional de gestiune: operațiile de modificare a BD sunt grupate în tranzații. Tranzațiile respecta **modelul ACID** (Atomica, Consistentă, Izolată, Durabilă)

SBD NoSQL : nu utilizează modelul ACID (datorită distribuției și replicării), și se înlocuiește cu **modelul BASE**



3. Infrastructura BD NoSQL

Modelul BASE:

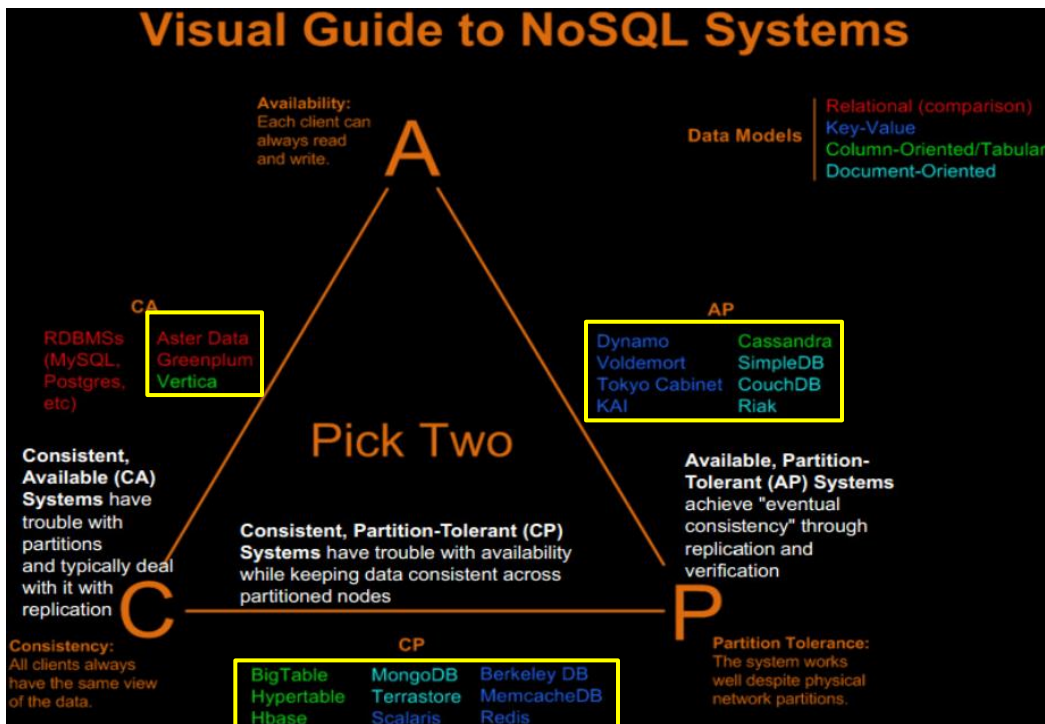
- ❑ **Basic Availability:** toți clienții primesc un răspuns la o interogare (nu dintr-o singură sursă a BD ci din colecția de date replicată și distribuită)
- ❑ **Soft State:** consistența bazei de date nu este verificată de SGBD, ea trebuie să fie asigurată de client/ aplicația care are dreptul de modificare a BD
- ❑ **Eventual Consistency:** BD poate fi inconsistentă (există valori diferite ale aceluiași date), dar în viitor datele vor ajunge într-o stare de consistență



4. Exemple SGBD NoSQL

Modelul	Exemple de sisteme
colecții de perechi cheie-valoare	Amazon Dynamo, Redis, Membase, MemcacheDB, Scalaris, Tokyo Cabinet, Voldemort, Riak
BigTable (column database)	Google Bigtable, Cassandra (Facebook), Hadoop/HBase, HyperTable, Amazon SimpleDB
graf	Neo4j, InfiniteGraph, InfoGrid, GraghBase, HyperGraphDB
colecții de documente	MongoDB, Couchbase, CouchDB, Terrastore, RavenDB

LIST OF NOSQL DATABASES [currently 225], <http://nosql-database.org/>





Clasamentul Top 10 al utilizării SGBD

310 systems in ranking, November 2016

Rank			DBMS	Database Model	Score		
Nov 2016	Oct 2016	Nov 2015			Nov 2016	Oct 2016	Nov 2015
1.	1.	1.	Oracle +	Relational DBMS	1413.01	-4.09	-67.94
2.	2.	2.	MySQL +	Relational DBMS	1373.56	+10.91	+86.71
3.	3.	3.	Microsoft SQL Server	Relational DBMS	1213.80	-0.38	+91.48
4.	↑ 5.	↑ 5.	PostgreSQL	Relational DBMS	325.82	+7.12	+40.13
5.	↓ 4.	↓ 4.	MongoDB +	Document store	325.48	+6.67	+20.87
6.	6.	6.	DB2	Relational DBMS	181.46	+0.90	-21.07
7.	7.	↑ 8.	Cassandra +	Wide column store	133.97	-1.09	+1.05
8.	8.	↓ 7.	Microsoft Access	Relational DBMS	125.97	+1.30	-14.99
9.	9.	↑ 10.	Redis	Key-value store	115.54	+6.00	+13.13
10.	10.	↓ 9.	SQLite	Relational DBMS	112.00	+3.43	+8.55

<http://db-engines.com/en/ranking>



Modele de BD NoSQL

- a) Colecții de perechi cheie-valoare
- b) Modelul BigTable
- c) Modelul graf
- d) Colectii de documente

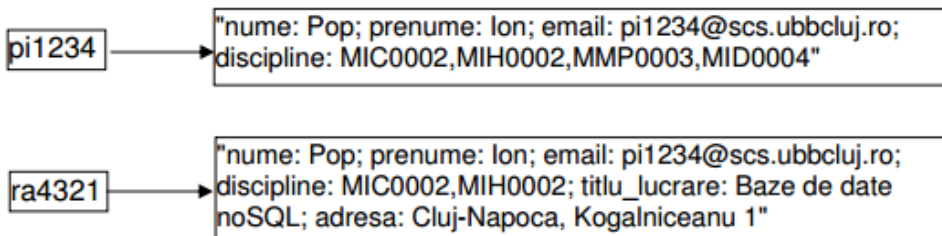


a) Colecții de perechi cheie-valoare

BD: colecție de perechi (cheie, valoare), cheia și valoarea =șiruri caractere, iar cheile sunt distincte (identificare).

Operații permise :

- adăugarea/ eliminarea unei perechi la/din colecție
- modificarea valorii dintr-o pereche existentă
- afișarea valorii pentru o cheie dată.



b) Modelul Big Table

Caracteristici:

- propus și folosit de Google** : Google Maps, Google Reader, Google Earth, Youtube și Gmail
- recomandat pentru tabele de dimensiuni mari**: multe elemente nedefinite (celule în tabel cu valori null), sau cu valori repetitive.
- la cerere /la modificarea valorii **se pot păstra alte versiuni ale valorilor** (timestamp).
- în BD relaționale valorile coloanelor memorate: linie, în BD BigTable **mod de memorare orientat coloană**. Există familii de coloane
- fiecare înregistrare se identifică prin valorile unei chei.



b) Modelul Big Table

Model orientat linie

ID	nume	prenume	email	contract_studiu
1	Pop	Ion	pi1234@scs.ubbcluj.ro	"MIC0002", "MIH0002", "MMP0003", "MID0004"
2	Popa	Radu		
3	Alb	Ana	aa4321@scs.ubbcluj.ro	"MLR0020", "MLR0002", "MLR5004"

Model orientat coloana: ID cheia tabelului

ID	nume	ID	prenume	ID	email	ID	contract_studiu
1	Pop	1	Ion	1	pi1234@scs.ubbcluj.ro	1	MIC0002
2	Popa	2	Radu	3	aa4321@scs.ubbcluj.ro	1	MIH0002
3	Alb	3	Ana			1	MMP0003
						1	MID0004
						3	MLR0020
						3	MLR0002
						3	MLR5004



b) Modelul Big Table

Daca se fac modificari in tabelul de mai jos

ID	nume	ID	prenume	ID	email	ID	contract_studiu
1	Pop	1	Ion	1	pi1234@scs.ubbcluj.ro	1	MIC0002
2	Popa	2	Radu	3	aa4321@scs.ubbcluj.ro	1	MIH0002
3	Alb	3	Ana			1	MMP0003
						1	MID0004
						3	MLR0020
						3	MLR0002
						3	MLR5004

Se utilizeaza timestamps: se introduce coloata ts = timp (ordine)modificare

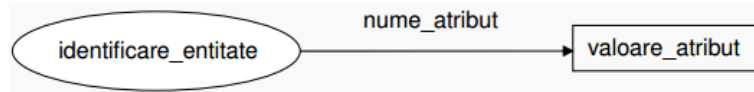
ID	ts	nume	ID	ts	email
1	t1	Pop	1	t1	pi1234@scs.ubbcluj.ro
2	t1	Popa	1	t2	piir1234@scs.ubbcluj.ro
3	t1	Alb	1	t3	pop_ion@yahoo.com
3	t5	Rus	3	t1	aa4321@scs.ubbcluj.ro
			3	t4	pop_ana@gmail.com



c) Modelul graf

Caracteristici:

- Bazat pe modelul graf: se construiesc o mulțime de triplete :
(identificare_entitate, nume_atribut, valoare_atribut)



Exemple

- Introducing the Knowledge Graph
<http://www.youtube.com/watch?v=mmQl6VGvX-c>
- The Knowledge Graph,
<http://www.google.com/insidesearch/features/search/knowledge.html>



c) Modelul graf

Exemplu: pentru BD cu tabelele de mai jos

studenti		
IDstud	nume	prenume
1	Pop	Ion
2	Popa	Radu
3	Alb	Ana

discipline	
cod	denumire
MIC0002	Sisteme de operare distribuite
MIH0002	Baze de date
MMP0003	Probabilități și statistică
MLR0002	Analiză matematică
MLR5004	Arhitectura sistemelor de calcul

contracte	
IDstud	coddisc
1	MIC0002
1	MIH0002
1	MMP0003
2	MIH0002
2	MLR5004
3	MLR0002
3	MLR5004

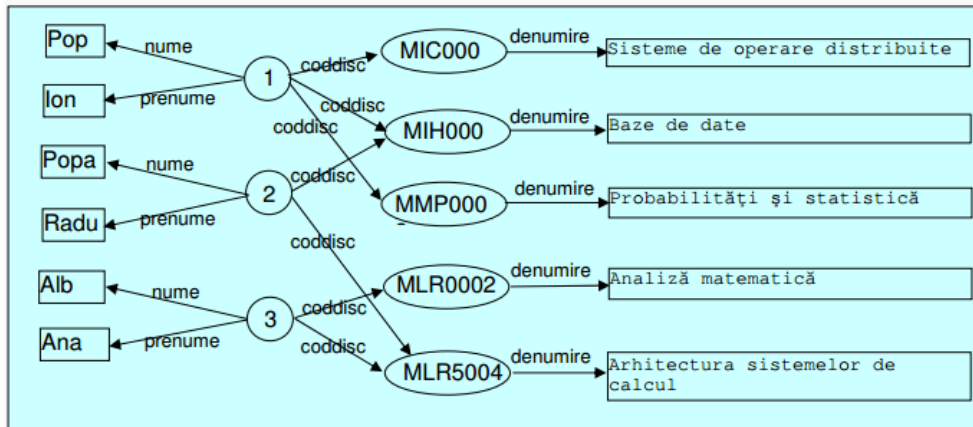
Tripletele sunt:

identificare_entitate	nume_atribut	valoare_atribut
1	nume	Pop
1	prenume	Ion
2	nume	Popa
2	prenume	Radu
3	nume	Alb
3	prenume	Ana
MIC0002	denumire	Sisteme de operare distribuite
MIH0002	denumire	Baze de date
MMP0003	denumire	Probabilități și statistică
MLR0002	denumire	Analiză matematică
MLR5004	denumire	Arhitectura sistemelor de calcul
1	coddisc	MIC0002
1	coddisc	MIH0002
1	coddisc	MMP0003
3	coddisc	MLR0002
3	coddisc	MLR5004



c) Modelul graf

Graful rezultat:



d) Colectii de documente

Colectie de documente: o BD care conține diverse colecții de documente /obiecte, analog tabelor dintr-o BD relațională.

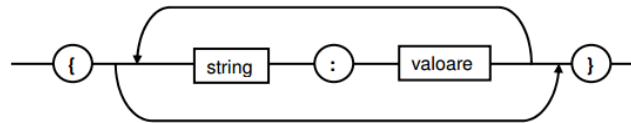
Caracteristici:

- într-o colecție se grupează documentele utile într-o interogare.
- între colecții diferite nu se pot efectua operații join.
- BD și colecțiile din BD se identifică printr-un nume (pentru numele bazei de date nu se poate folosi: admin, local, config).
- un document (obiect) nu are o structură stabilită
- un document are o identificare unică în colecția de date, prin câmpul cu denumirea "_id".
- caracterele (din denumiri, din valori) sunt case-sensitive
- pentru gestiunea unei BD și a unei colecții există mai multe metode



d) Colectii de documente

Document: format dintr-o mulțime de perechi (câmpuri) nume/valoare:



Această structură de document se folosește pentru:

- documentele memorate în baza de date
- precizarea condițiilor pe care le îndeplinesc documentele utile într-o comandă (citire, actualizare, ștergere) { string : valoare }
- precizarea tipurilor de modificări pentru documente
- specificarea modului de construire a indexurilor
- diverse opțiuni în comenzi



d) Colectii de documente

Exemplu: declarare document in BD

```
{
  "nume": "Pop",
  "prenume": "Ion",
  "contract_studiu": ["MIC0002", "MIH0002", "MMP0003", "MID0004"],
  "adresa": {
    "localitatea": "Cluj-Napoca",
    "strada": "Kogalniceanu",
    "numarul": 1
  }
  "email": "pi1234@scs.ubbcluj.ro"
}
```



Concluzii

Avantaje utilizare BD NoSQL:

- Nu sunt un substitut al BD relationale, dar sunt dedicate rezolvarii unor probleme specifice ale acestora
- indicate cel mai frecvent: pentru cresterea eficientei si disponibilitatii in detrimentul consistentei BD
- Utilizeaza modelul BASE nu ACID
- Utile pentru BD de dimensiuni mari
- Modelul de date este flexibil
- Scalabilitate orizontala
- Majoritatea sunt free/open source
- Usor de instalat si configurat



Concluzii

Dezavantaje:

- lipsa standardelor
- tehnologii in curs de dezvoltare
- dificil de administrat
- putini experti in domeniu