

# BIG DATA – CONCEPTE, ARHITECTURI ȘI TEHNOLOGII

Adriana ALEXANDRU

adriana@ici.ro

Dora COARDOȘ

coardos@ici.ro

Institutul Național de Cercetare-Dezvoltare în Informatică - ICI București

**Rezumat:** Lucrarea prezintă principalele abordări legate de Big Data. Conceptul Big Data introduce modificări în cadrul a trei dimensiuni: (1) tipuri de date, (2) viteza de acumulare a acestora și (3) volumul lor. Big Data reprezintă o nouă generație de tehnologii și arhitecturi destinate extragerii de valoare din cadrul volumelor foarte mari de date, cu o mare varietate, permițând prelucrarea și analiza acestora în timp real. În articol sunt prezentate seturi de tip Big Data, precum și arhitecturi și tehnologii pentru Big Data.

**Cuvinte cheie:** Big Data, seturi Big Data, Hadoop, NoSQL.

**Abstract:** This paper presents the main approaches related to Big Data. Big Data concept introduces changes in three dimensions: (1) data types, (2) speed of data accumulation and (3) data volume. Big Data represents a new generation of technologies and architectures for extracting value from the very large volumes of data that have a wide variety, and allow real time processing and analysis. This article presents Big Data sets, architectures, and technologies.

**Keywords:** Big Data, Hadoop, NoSQL, Big Data sets.

## 1. Introducere

În ultimele decenii, organizațiile au început să acorde importanță sporită datelor și să investească mai mult în colectarea și gestionarea lor. Dincolo de informațiile colectate în interiorul organizațiilor și de volumul crescând de date pe care le generează calculatoarele în funcționarea lor, sunt utilizate date obținute din exteriorul organizației, fie structurate sau nestructurate, care au surse multiple care pot include de la informații postate pe rețele de socializare și produse vizionate în magazine virtuale, la informații citite de către senzori, semnale GPS de pe dispozitivele mobile, adrese IP ale compu-terelor, cookie-uri, coduri de bare ș.a.m.d. Unele tipuri de date precum text și voce, există de mult timp, însă volumul acestora în mediul Internet și în alte structuri digitale anunță începutul unei noi ere, precum și a unor noi tehnologii care permit analizarea acestor tipuri de date.

Multe dintre cele mai importante surse de date sunt însă relativ noi. Se argumentează că explozia volumului de date caracteristic fenomenului Big Data provine din datele de natură nestructurată. În cadrul acestora, spre deosebire de datele generate de către utilizatori, care au la origine informații furnizate voluntar în diferite medii de diseminare Web, există și datele interceptate. Acestea din urmă se referă la informații colectate în mod pasiv din comportamentul online al indivizilor, cum sunt, de pildă, termenii de căutare online sau

localizarea indivizilor prin aplicațiile prezente pe dispozitivele mobile.

În comparație cu instrumentele analitice tradiționale, conceptul Big Data introduce modificări în cadrul a patru dimensiuni: (1) tipuri de date, (2) viteza de acumulare a acestora, (3) volumul lor și (4) calitatea datelor. Odată cu lansarea mediului Web 2.0, o mare parte din datele de valoare pentru întreprinderi sunt generate în exteriorul organizației, de către consumatori și în general, utilizatori Web.

Ne aflăm astăzi la un punct de inflexiune în care volumul și varietatea datelor generate în organizații creează provocări, dar și oportunitatea de a atinge noi posibilități de afaceri și valoare adăugată. Cei care vor avea capacitatea de a construi infrastructurile corespunzătoare pentru managementul informației vor putea transforma aceste provocări în avantaj competitiv și vor putea propulsa afacerea lor spre rezultate mult mai bune. Adăugarea de valoare prin utilizarea potențialului Big Data este încă în faza emergentă, însă reprezintă o schimbare de paradigmă pe care orice afacere trebuie să o ia în considerare.

Approape orice companie descoperă că trebuie nu doar să gestioneze volume de date din ce în ce mai mari în sistemele lor în timp real, dar și să analizeze aceste informații astfel încât să poată lua rapid deciziile potrivite pentru a concura eficient pe piață.

Cererea crescândă pentru platforme

analitice de generație următoare care oferă clienților răspunsuri aproape în timp real, declanșate de date în timp real cum ar fi istoria accesărilor sau parcursul vizitatorilor (clickstreams), social media, senzori, combinate cu puterea de executare distribuită a seriilor de comenzi, demonstrează faptul că inteligența ar trebui să se afle, implicit, în centrul oricărei aplicații software. Astfel, aplicațiile moderne pun accent pe utilizarea noilor tehnologii Big Data.

Soluția oferită prin utilizarea seturilor Big Data ajută afacerile să își administreze mai bine fluxurile de date cu volum mare, varietate mare și viteză mare și să transforme aceste date în informații care să genereze profit.

## 2. Conceptul Big Data

### 2.1 Definiții și caracteristici

În fiecare zi, creăm un număr mare de date. Aceste date sunt de peste tot: senzori care înregistrează vremea, post-uri pe social media, fotografiile digitale și clipuri video, tranzacții online, semnale GPS de pe telefon și multe altele. Acestea reprezintă ”Big Data”.

*Conceptul de Big Data* este în prim-planul temelor actuale în cele mai multe cercuri de IT. Înțelegerea conceptului de Big Data, la fel ca orice altă tehnologie în curs de dezvoltare, necesită mai întâi ca acesta să fie definit.

Conform [1], există cel puțin 43 definiții ale termenului Big Data. Câteva sunt prezentate în continuare:

- în 2011, un raport al International Data Corporation a definit Big Data ca fiind “o nouă generație de tehnologii și arhitecturi, proiectate pentru a extrage valoare economică din volume foarte mari de date de o largă varietate, prin asigurarea unei viteze ridicate de captare, descoperire și/sau analiză” [2];
- Gartner definește Big Data fiind "date de volum mare, de mare viteză și de mare varietate care necesită forme inovatoare rentabile de prelucrare a informațiilor pentru a intensifica înțelegerea domeniului și luarea deciziilor" [3].

Cercetătorii domeniului Big Data sunt în unanimitate de acord că toate sistemele de tipul Big Data au următoarele caracteristici

definitorii pentru datele lor: *volumul (volume)*, *varietatea (variety)*, *viteza (velocity/virality)*, *veridicitatea (veracity)*, *validitatea (validity)*, *variabilitatea (variability)*, *volatilitatea (volatility)*, *vâscozitatea (viscosity)*, *vizualizarea (visualization)*, și *valoarea (value)*. Trebuie să menționăm că numai primii "patru V", dacă au valori “mari”, definesc Big Data, restul celor "șase V" se regăsesc la orice fel de date.

Întrucât cei "patru V" sunt considerați definitorii pentru acest concept, este oportună o detaliere a semnificației acestor caracteristici:

1. *Volum*: creșterea volumelor de date în sisteme de tip întreprindere este cauzată de volumul tranzacțiilor și a altor tipuri de date tradiționale, precum și de noi tipuri de date. Un volum prea mare de date reprezintă o problemă de stocare, dar prea multe date au în egală măsură și un mare impact asupra complexității analizei datelor;

2. *Viteză*: se referă atât la rapiditatea cu care datele sunt produse, cât și la rapiditatea cu care datele trebuie să fie prelucrate pentru a satisface cererea. Acest lucru implică fluxuri de date, crearea de înregistrări structurate, precum și disponibilitatea pentru acces și livrare. Viteza de generare, prelucrare și analiză a datelor crește continuu, în principal din următoarele motive: specificul de timp real al proceselor de generare, cererile care rezultă din combinarea fluxurilor de date cu procesele de afaceri, specificul proceselor de luare a deciziilor. Viteza de prelucrare a datelor trebuie să fie ridicată, în timp ce capacitatea de prelucrare depinde preponderent de tipul de prelucrare al fluxurilor de date;

3. *Varietate*: liderii IT au avut întotdeauna o problemă cu transformarea volumelor mari de informații tranzacționale în decizii, deși tipurile de date generate sau prelucrate erau puțin diversificate, mai simple și majoritar structurate. În prezent, există mai multe tipuri de informații pentru analiză generate de noile canale și tehnologii apărute - în principal provenind din *social media*, *Internetul lucrurilor*, *surse mobile* (sensibile la context) și publicitatea online – care generează date semistructurate sau nestructurate. Varietatea include date tabelare (baze de date), date ierarhice, documente, XML, e-mailuri, blog-uri, mesaje instant, click stream-uri, fișiere log, date de contorizare, imagini statice, audio,

video, date despre cursul acțiunilor (stoc ticker), tranzacții financiare etc.

4. *Veridicitate*: se referă la cât de încredere sau de îndoielnice sunt datele. Calitatea datelor Big Data este mai puțin controlabilă deoarece provine din diferite surse pentru care nu se poate garanta calitatea conținutului și forma lui de prezentare. Pentru analistul de date experimentat, este esențială capacitatea de a evalua conformitatea, acuratețea și sinceritatea datelor supuse analizei. Aici discuția se poartă în jurul responsabilității generatorului inițial al datelor, scopului pentru care datele sunt emise și reacțiilor receptorilor.

În anticiparea oportunităților Big Data, companiile din toate mediile industriale colectează și stochează provizoriu un număr imens de date operaționale, publice, comerciale sau sociale. În majoritatea mediilor, în special guvernamentale, producție și educație, combinarea acestor surse cu "dark data", cum ar fi email-uri, multimedia etc., reprezintă de cele mai multe ori cea mai nouă oportunitate de a transforma afacerile.

## 2.2 Seturi Big Data

Pentru procesarea Big Data, datorită datelor de complexitate și dimensiune foarte mare, nu pot fi utilizate aplicații standard fiind necesare aplicații capabile să ruleze în mod paralel pe un număr foarte mare de servere. Printre dificultățile întâlnite în procesarea acestor date se numără: capturarea, curățarea, stocarea, căutarea, partajarea, transferul, analiza și vizualizarea.

Sistemele de Big Data pot furniza informație atât organizațiilor guvernamentale cât și cetățenilor, provenind din diferite surse care pot fi identificate după cum urmează: document pe hârtie (mediu fizic), documente digitale, puncte de acces la rețeaua de Internet guvernamentală, site-uri localizate pe platformele online de socializare și sisteme operaționale disponibile.

Informația furnizată de sistemele Big Data nu include informații personale sau informații restricționate de mecanisme de control și confidențialitate.

Potrivit Garter [4], pentru a gestiona un volum mare de date, informațiile ar putea fi incluse în categorii, în funcție de sursă. Firma de consultanță americană a identificat cinci

astfel de tipuri de informații:

1. *Date operaționale*: sunt date despre consumatori, furnizori, parteneri și angajați deja accesibile pe baza unor procese de tranzacție sau din baze de date;

2. *Date ascunse (Dark Data)*: sunt informațiile adunate de-a lungul vremii în arhive, dar care nu pot fi clar structurate. Ele pot fi utilizate ulterior pentru luarea de decizii, analize de afaceri, etc. În acest caz ar fi incluse mail-urile, contractele, informațiile multimedia;

3. *Date comerciale*: sunt date care pot veni prin intermediul agregatoarelor de date (care citesc RSS-urile) specifice, în funcție de industrie;

4. *Date publice*: sunt datele publice care aparțin instituțiilor statului (informații care vin de la Guvern, de la ministere);

5. *Date din social media*: sunt datele care arată activitatea unui utilizator pe un blog, pe rețelele de socializare. Ele sunt utile pentru a stabili trenduri, atitudini, preferințe.

Big Data reprezintă seturi mari de informații complexe care în urma unei analize pot determina creșterea inteligenței în afaceri prin identificarea trendurilor și îmbunătățirea operațiunilor de afaceri și proceselor decizionale, pot contribui la prevenirea bolilor și chiar combate rata criminalității.

Dintre domeniile în care proiectele Big Data sunt realizabile amintim: Sănătate (analiza statistică a cazurilor, telemedicină etc.), Cultură, eCommerce, Securitate națională.

## 3. Arhitecturi și tehnologii Big Data

### 3.1 Arhitecturi pentru sistemele Big Data

Din cauza complexității sistemelor Big Data, a fost necesară dezvoltarea unei *arhitecturi specializate*. BDAF (Big Data Architecture Framework) are ca scop implementarea unei colecții specifice de elemente de design, asigurarea abordării unui design consistent, reducerea complexității sistemului, maximizarea reutilizării, legăturilor slabe (loose-coupling), reducerea dependențelor și creșterea productivității.

Pentru Big Data, cea mai frecventă arhitectură utilizată este Hadoop. Această inovație a redefinit managementul datelor, deoarece prelucrează cantități mari de date, cu costuri reduse și în timp util.

### 3.1.1 Ecosistemul Hadoop

Fiind o colecție de seturi de date atât de mare și de complexă, utilizarea uneltelor manuale de gestionare a bazelor de date devine incomodă în cazul Big Data. Atunci când lucrăm cu volume mari de date avem nevoie de o soluție care să ne permită atât stocarea la un cost cât mai mic, dar și să asigure o performanță bună la procesare. Un posibil răspuns la această provocare este platforma de aplicații *Apache Hadoop*.

*Hadoop* [5] este un proiect open-source dezvoltat de Apache care își propune realizarea de procesări distribuite a unor seturi de date de dimensiuni mari, rulând pe mai multe clustere, folosind modele de programare simple. Proiectarea acestui framework a fost realizată astfel încât să fie scalabilă chiar și în situația în care sarcinile sunt rulate pe mii de calculatoare, fiecare dintre acestea punând la dispoziție o anumită capacitate de procesare și de stocare.

Începând cu anul 2010, Hadoop a fost adoptat pe scară largă de organizații atât în scopul de a stoca volume mari de date, cât și ca platformă de analiză a acestora. În prezent, Hadoop este folosit de numeroase companii pentru care volumul de date generat zilnic depășește capacitățile de procesare și stocare specifice sistemelor convenționale: Adobe, AOL, Amazon.com, EBay, Facebook, Google, LinkedIn, Twitter, Yahoo.

Apache Hadoop este un *ecosistem* de unelte gândite pentru a funcționa împreună ca o soluție eficientă de stocare și procesare a datelor. Aceste unelte sunt dezvoltate de către o comunitate diversificată de dezvoltatori într-un mod colaborativ sub umbrela Apache Software Foundation.

Produsele Hadoop integrate în cele mai multe dintre distribuții sunt HDFS, MapReduce, HBase, Hive, Mahout, Oozie, Pig, Sqoop, Whirr, ZooKeeper, Flume [6].

Nucleul Apache Hadoop este format din **două componente**: un sistem de fișiere distribuit (*HDFS – Hadoop Distributed File System*) și un framework pentru procesare

distribuită (*MapReduce*). Hadoop a fost gândit să funcționeze într-o arhitectură de tip cluster construită pe echipamente server obișnuite.

- **HDFS** oferă o stocare extrem de fiabilă și distribuită, prin replicarea datelor pe mai multe noduri. Spre deosebire de un sistem de fișiere obișnuit, atunci când datele sunt trimise la HDFS, acestea se vor împărți în mod automat în mai multe blocuri și depozitează datele în diferite „DataNodes”. Acest lucru asigură disponibilitate ridicată și toleranță la erori.
- **MapReduce** oferă un sistem de analiză care poate efectua calcule complexe, pe seturi de date de dimensiuni mari. O procesare de tip MapReduce presupune că problema care trebuie rezolvată poate să fie împărțită în probleme mai mici care pot să fie rezolvate independent (faza de *map*), într-o manieră “divide et impera”, fiecare fiind executată cât mai aproape de datele pe care trebuie să opereze urmând ca apoi rezultatele să fie reunite în funcție de necesități (faza de *reduce*) [7]. Această componentă este responsabilă de efectuarea calculelor și de împărțirea unui calcul de complexitate ridicată în mai multe task-uri.
- **HBase** este o bază de date distribuită de tip NoSQL, orientată pe coloane având la baza modelul Google BigTable, care folosește ca și mediu de stocare HDFS, fiind utilizată în cazul aplicațiilor Hadoop care necesită operații de citire / scriere aleatoare în seturi de date foarte mari. HBase reprezintă o soluție pentru seturi de informații de dimensiuni foarte mari (de ordinul milioane și miliardelor de înregistrări) sau pentru aplicații ce utilizează date care sunt accesate de foarte mulți clienți (cererile și răspunsurile generate ca urmare a acestei interacțiuni implică un volum de date foarte mare). Totodată, funcționează optim în cazul unor scheme variabile, unde structura înregistrărilor diferă (datorită unor atribute care pot să existe sau nu). HBase are trei componente principale: biblioteca clientului, un server de tip master, mai multe servere de regiune.
- **Zookeeper** este un serviciu de coordonare pentru aplicațiile distribuite. Zookeeper menține, configurează și denumește

cantități mari de date. De asemenea, furnizează servicii distribuite de sincronizare și de grup. Singur, Zookeeper conține noduri master și slave și stochează informații de configurare.

- **HCatalog** stochează metadate și generează tabele pentru cantități mari de date. HCatalog simplifică comunicarea utilizator folosind datele HDFS și este o sursă de partajare a datelor între instrumente și platformele de execuție.
  - **Hive** este o platformă de depozitarea datelor (de tip data warehouse) care permite interogarea și gestionarea seturilor de date de mari dimensiuni din depozite distribuite, stocate în HDFS. Hive este o sub-platformă în ecosistemul Hadoop și folosește un limbaj de interogare de tipul SQL, care este numit HiveQL. Limbajul, de asemenea, permite programatorilor tradiționali ai MapReduce să se conecteze la mediul lor specific de interogare și de reducere atunci când este incomod sau ineficient.
  - **Pig** este o platformă de nivel înalt folosită pentru analizarea unor seturi de date mari având un limbaj propriu, pentru descrierea programelor de analiză a datelor. Caracteristica principală a Pig este că prin natura programelor Pig, permite paralelizarea lor la momentul rulării. Compilatorul Pig produce joburi MapReduce. Arhitectura Pig generează un limbaj de scripting de nivel înalt (Pig Latin) și operează pe o platformă în timp real, platformă care permite utilizatorilor să execute MapReduce pe Hadoop. Pig este mai flexibil decât Hive referitor la formatul datelor, furnizând propriul model de date. Pig are propriul tip de date, hartă, care reprezintă datele semistructurate, inclusiv JSON și XML.
  - **Mahout** este o bibliotecă pentru algoritmi de învățare automată (machine-learning) și data mining, incluzând algoritmi de clasificare și de clustering. Mulți algoritmi sunt scriși pentru compatibilitate cu MapReduce, astfel încât ei sunt scalabili la seturi de date mari. Această componentă este împărțită în patru grupe principale: filtrare colectivă, clasificare, clustering și extragere de modele paralele frecvente (mining of parallel frequent patterns).
  - **Oozie** este un instrument pentru managementul workflow-ului /coordonarea joburilor MapReduce. Apache Oozie permite combinarea mai multor elemente într-o unitate logică de lucru. Apache Oozie este o aplicație Java Web, care rulează într-un servlet Java Tomcat și folosește o bază de date pentru a stoca definiții ale fluxului de lucru și execuții curente ale fluxului de lucru.
  - **Avro** serializează datele, conduce apelurile de proceduri la distanță și transferă datele de la un program sau limbaj la altul. În această arhitectură, datele se auto-descriu și sunt întotdeauna stocate în funcție de propria lor schemă, deoarece aceste calități sunt potrivite limbajelor de scripting, cum ar fi Pig.
  - **Chukwa** este un tool pentru monitorizarea aplicațiilor distribuite, bazându-se pe arhitectura HDFS și Map Reduce. Chukwa este o arhitectură pentru colectarea și analiza datelor. Chukwa colectează și prelucrează datele din sistemele distribuite și le stochează în Hadoop.
  - **Flume** este un serviciu distribuit care permite colectarea, agregarea și mutarea unor volume mari de date tip log. Are o arhitectură bazată pe fluxuri de date și care permite construirea de aplicații analitice. Folosește două canale, și anume, surse și colectoare (sinks). Sursele includ date Avro, fișiere și fișierele jurnal (log) de sistem, în timp ce colectoarele fac referire la HDFS și HBase. Prin motorul său personal de prelucrare, interogare, Flume transformă fiecare nou batch de Big Data înainte de a fi transportat în colector.
- Deși Hadoop are diverse componente, fiecare companie utilizează anumite componente ale Hadoop în funcție de necesitățile ei.

### 3.1.2 Integrare Big Data cu Hadoop

Arhitectura de Big Data nu este una fixă, care să se potrivească în toate situațiile. Fiecare strat de procesare în arhitectură are mai multe soluții și tehnici care pot fi implementate pentru a crea un mediu robust. Fiecare soluție are propriile avantaje și dezavantaje pentru un anumit volum de muncă [8].

În general orice arhitectură de date se compune din patru componente logice principale: sursele de date, transformarea datelor, prelucrarea datelor sau integrarea datelor, cereri de date. Printre problemele legate de integrarea Big Data se numără varietatea surselor de date, calitatea datelor ce urmează a fi integrate și vizualizarea datelor [9]. Arhitectura unui ecosistem pentru integrarea Big Data (vezi figura 1) include următoarele componente:

1. *Sursele de date structurate și nestructurate.* Introducerea bazelor de date stocate în cloud și a infrastructurii mobile, au dus la o creștere semnificativă a dimensiunii și complexității seturilor de date, acestea devenind componenta principală a ecosistemelor de integrare a datelor. Astfel arhitectura de integrare a datelor trebuie să includă strategii multiple pentru accesarea și stocarea unei cantități foarte mari și diversificate de date.

stocarea datelor și să realizeze conexiunile cu celelalte surse de date.

2. *Platforma pentru descoperirea datelor.* Platforma pentru descoperirea datelor reprezintă un set de instrumente și tehnici pentru lucrul cu fișiere pentru Big Data pentru găsirea de modele și răspunsuri la întrebări de business. În prezent, aceasta este mai mult o activitate adhoc, iar organizațiile întâmpină dificultăți în dezvoltarea unor procese în jurul ei. În cadrul activității de descoperire a datelor, informațiile obținute pot deveni uneori neutilizabile după doar câteva ore. Arhitectura pentru integrarea datelor trebuie să țină cont de aceste informații volatile pentru asigurarea calității datelor. Infrastructura pentru integrarea datelor trebuie să fie capabilă să răspundă rapid la cerințele utilizatorilor.
3. *Depozite de date tradiționale.* Depozitele de date tradiționale oferă necesarul de

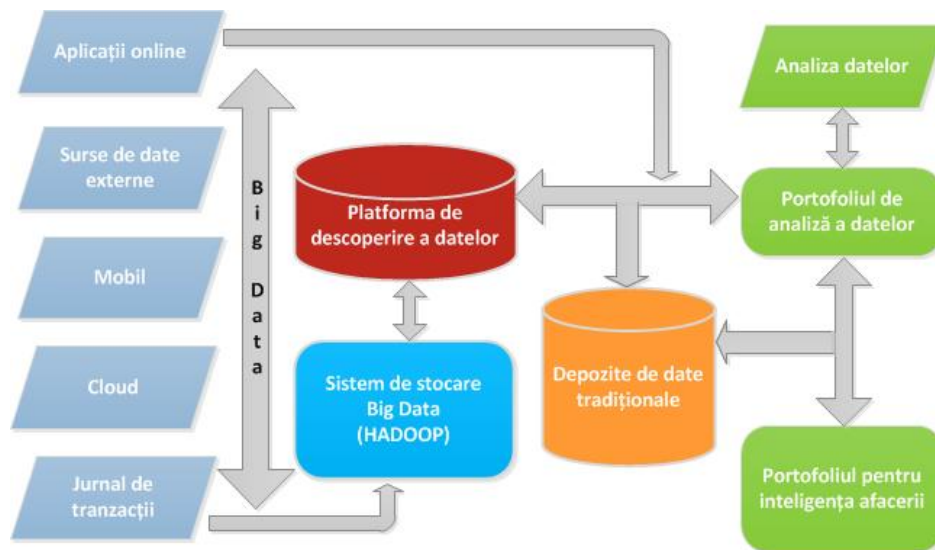


Figura 1. Arhitectura unui ecosistem pentru integrarea Big Data [9]

*Sisteme de stocare Big Data.* În timp ce sistemele de stocare a datelor foarte mari precum Hadoop asigură mijloace de stocare și organizare a unor volume mari de date, procesarea acestora pentru extragerea de informații utile rămâne în continuare o activitate dificilă. Arhitectura MapReduce a acestor sisteme dă posibilitatea de stocare rapidă a unor cantități foarte mari de date și oferă suport pentru realizarea de analize pe baza acestor date. Platforma pentru integrarea datelor trebuie să construiască structura pentru

informații de bază, dar trebuie să includă noi funcționalități pentru o mai bună integrare a surselor de date nestructurate și pentru a satisface nivelul de performanțe solicitat de platformele de analiză. Organizațiile au început să dezvolte noi modalități de separare a analizelor operaționale de analizele în profunzime pe baza istoricului pentru deciziile strategice. Platforma pentru integrarea datelor trebuie să fie capabilă să separe informațiile operaționale, de sursele de date utilizate în elaborarea strategiilor pe termen lung. Totodată infrastructura de

integrare a datelor trebuie să permită un acces rapid la datele cel mai des accesate.

4. *Portofoliul pentru inteligența afacerii.* Portofoliul pentru inteligența afacerii se concentrează pe rezultatele și performanțele din trecut, chiar dacă va exista o creștere a cererii pentru rapoarte și performanțe operaționale. Evoluția necesității de autoservire a inteligenței afacerii și inteligenței afacerii pe dispozitive mobile va continua să genereze probleme arhitecturale platformelor de integrare a datelor. Un alt aspect foarte important îl reprezintă capacitatea portofoliului de inteligență a afacerii de integrare cu portofoliul de analiză. Aceasta poate conduce la o creștere a cererilor pentru integrarea informațiilor.
5. *Portofoliul de analiză a datelor.* Activitatea de analiză din cadrul acestui portofoliu trebuie să gestioneze atât problemele legate de activitatea companiei, cât și cele legate de date. Platformele de integrare a datelor joacă două roluri în ceea ce privește asigurarea suportului necesar portofoliului de analiză. În primul rând, ecosistemul de integrare a datelor trebuie să asigure accesul la date structurate și nestructurate pentru activitatea de analiză. În al doilea rând, trebuie să permită reutilizarea analizelor efectuate anterior, reducând astfel situațiile care ar necesita repetarea unor pași.

Ecosistemul de integrare a datelor trebuie deci să includă posibilitatea de procesare a unor volume foarte mari de date și să facă față unor solicitări de a lucra cu o varietate mare de surse de date.

Una dintre principalele provocări, în ceea ce privește prelucrarea de seturi foarte mari de date, este manipularea fluxurilor de date în timp real. În timp ce ambele tipuri de date, offline și online, pot fi în mod independent prelucrate, adesea este nevoie să furnizăm răspunsuri la întrebările cu privire la evenimente online bazate pe trecut. **Arhitectura Lambda** vine ca un răspuns la aceste provocări [10]. Această arhitectură integrată cu Hadoop este prezentată în figura 2.

Arhitectura, așa cum este prezentată în figura 2, este alcătuită din următoarele componente:

- *Stratul de loturi* (Batch Layer) - responsabil pentru gestionarea setului de date master și de precalcularea vizualizărilor batch;
- *Stratul de servire* (Serving Layer) - indexează vizualizările batch pentru interogări ad-hoc;
- *Stratul de viteză* (Speed Layer) - servește doar datelor noi, care nu au fost încă procesate de Nivelul batch.

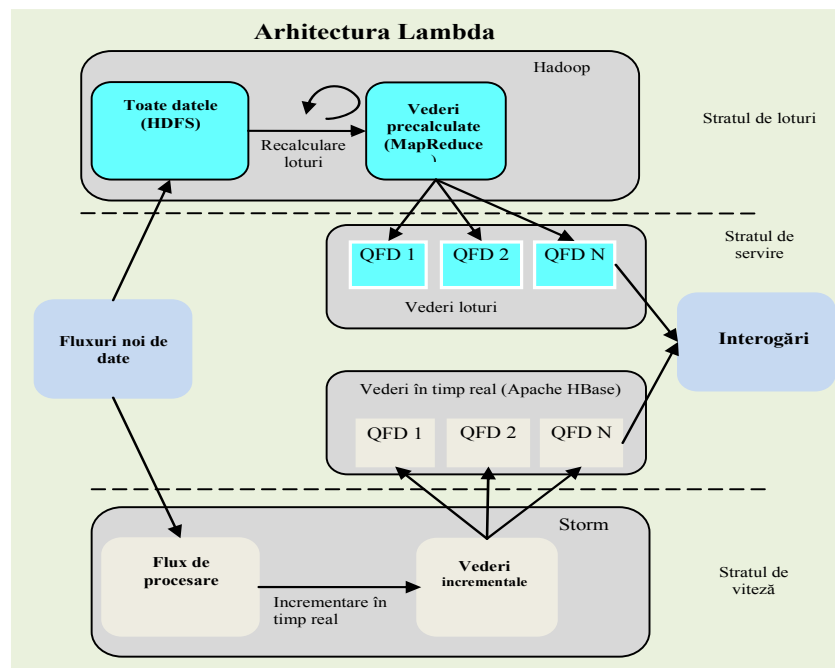


Figura 2. Arhitectura Lambda integrată cu Hadoop [10]

## 3.2 Tehnologii pentru Big Data

*Tehnologiile Big Data* reprezintă un domeniu aflat în continuă dezvoltare, ce se ocupă cu analiza și gestionarea volumelor mari de date. Această definiție cuprinde atât echipamentele hardware, cât și sistemele software care integrează, organizează, gestionează, analizează și prezintă Big Data.

Unele tipuri de date precum text și voce, există de mult timp, însă volumul acestora în mediul Internet și în alte structuri digitale anunță începutul unei noi ere, precum și a unor noi tehnologii care permit analizarea acestor tipuri de date.

Tehnologiile Big Data pot fi clasificate în șase categorii [11]):

### 1. Tehnologii suport pentru infrastructură

constau în:

- *platformele Cloud pentru Big Data* – au resurse eficiente de programare și gestionare,
- *tehnologii de stocare* - legate de compresia datelor și de virtualizare de stocare,
- *tehnologii de virtualizare* – se referă la procesul de partajare a resurselor și de izolarea hardware-ului de la bază,
- *tehnologii de rețea* - măresc considerabil frecvența și viteza de transmisie a datelor,
- *tehnologii de monitorizare a resurselor* - gestionează resursele conectate la rețea în scopul identificării erorilor apărute în sistem;

**2. Tehnologii pentru achiziționarea datelor** - obținerea datelor neprelucrate (brute) de la senzori sau alte surse dedicate prin intermediul:

- *tehnologiei de achiziționare a datelor bazată pe senzori* - care permite ca informațiile oferite de senzori să fie transferate către o bază de date cu ajutorul rețelelor wireless,
- *tehnologiei de achiziționare a datelor bazată pe rețele de date;*

**3. Tehnologii pentru transferul datelor** - pentru colectarea informațiilor înainte de procesarea datelor;

**4. Tehnologii pentru memorarea și arhivarea datelor** - se aplică în mod distribuit

în noduri de stocare multiple, pentru care se pun la dispoziție mecanisme de back-up, securitate, interfețe de acces și protocoale și includ următoarele:

- *sisteme de fișiere distribuite:* pentru prelucrarea datelor trebuie adoptată, de asemenea, o arhitectură și soluții distribuite – *HDFS, HBase, Cassandra, MongoDB* open source,
- *baze de date relaționale:* tradiționale, caracterizate prin lipsă de scalabilitate și extensibilitate, nu sunt adecvate. Interogările pe bazele de date SQL (*MySQL* și *Oracle*) de pe disc sunt lente,
- *tehnologii NOSQL:* noi tehnologii pentru baze de date ce nerelaționale, care nu oferă garanțiile ACID (Atomicitate, Consistență, Izolare) și sunt utilizate în prelucrarea datelor nestructurate și analiza Big Data;

### 5. Tehnologii pentru procesarea datelor

- se referă la aspectele de procesare a datelor și utilizarea tehnicilor de bază ale tehnologiilor Big Data, pentru analizarea, prelucrarea și exploatarea datelor, extragerea de informații și cunoștințe importante și apoi transformarea în modele utile și aplicarea acestora la procesele de cercetare și operare;

**6. Tehnologii pentru afișarea datelor și interacțiune** - urnizarea de vizualizări interactive, care permit utilizatorilor să navigheze prin seturile de date. Permite utilizatorilor să ia decizii care să sprijine producția, operarea și planificarea.

### 3.2.1 Tehnologia NoSQL

Bazele de date relaționale tradiționale nu pot face însă față provocărilor actuale aduse de către Big Data. În ultima vreme bazele de date de tipul NoSQL sunt din ce în ce mai populare pentru stocarea datelor de mari dimensiuni. Au apărut din necesitatea unor companii precum Google, Facebook sau Twitter de a manipula cantități imense de date cărora bazele de date tradiționale pur și simplu nu le pot face față. Așa că bazele de date NoSQL au fost proiectate pentru a stoca volume foarte mari de date în general fără o schemă fixă și partiționate pe multiple servere. Bazele de date NoSQL oferă moduri flexibile de lucru, suport pentru copierea datelor mult mai simplu și mai ușor, un API simplu, și coerența eventuală a datelor. Bazele de date NoSQL devin astfel



tehnologia de bază pentru Big Data.

NoSQL (Not Only SQL) sunt baze de date non relaționale [12]. Principalul avantaj al utilizării bazelor de date NoSQL este acela că permit lucrul eficient cu date structurate, precum e-mailul, multimedia, procesoare de text. Bazele de date NoSQL, ca nouă generație de baze de date: nu sunt relaționale, sunt distribuite, sunt Open Source și se caracterizează prin scalabilitate orizontală. O altă caracteristică importantă a sistemelor NoSQL este arhitectura "shared nothing" prin care fiecare nod-server este independent, nu partajează memorie sau spațiu.

Bazele de date NoSQL au o structură mai simplă și o tehnologie diferită pentru stocarea și extragerea datelor decât bazele de date relaționale și oferă performanțe mai bune pentru analize în timp real sau pe volume mari de date. Într-o bază de date NoSQL nu există o schemă propriu-zisă a datelor, ele fiind stocate ca perechi cheie-valoare (foarte eficient și flexibil, dar datele nu sunt self-describing), sau de coloane (folosit pentru date împrăștiate), sau document (folosit pentru depozite XML, dar ineficient ca performanță), sau graf (folosit pentru traversări relaționate, dar ineficient la căutări) [13].

Astfel mișcarea NoSQL reprezintă o încercare de a depăși limitările modelului relațional și un pas de trecere către NewSQL și anume relațional plus extra funcționalități NoSQL.

Cele mai populare baze de date NoSQL în acest moment sunt: Cassandra, MongoDB, CouchDB, Redis, Riak, Membase, Neo4j și HBase.

#### 4. Concluzii

Big Data poate adăuga valoare și oferi o nouă perspectivă prin îmbunătățirea practicilor de analiză și modelare predictivă. Volumele masive de date provenind din surse diferite au un efect pozitiv în procesul de luare a deciziilor în timp real. Varietatea surselor de date, calitatea datelor care urmează să fie integrate și vizualizarea lor sunt unele dintre provocările pentru integrarea Big Data.

Noi capacități și tehnologii trebuie să fie adoptate în scopul de a transforma informațiile prin gestionarea și de analiza datelor. Principalele provocări sunt acceptarea și

utilizarea noilor tehnologii, precum și reglementarea lor. Cele mai notabile probleme de depășit rezidă în dificultatea de a analiza volume mari de date pentru a obține rezultate precise în timp util, necesitatea de standardizare, interoperabilitatea, securitatea, confidențialitatea, precum și expertiza și finanțarea pentru dezvoltarea infrastructurii Big Data și integrarea seturilor de date deja disponibile.

Noile metode, instrumente și abordări statistice tehnologice, cum ar fi cloud computing și tehnologii de securitate trebuie să fie explorate. Mai mult decât atât, ar trebui să se investească în instruirea personalului cu privire la utilizarea BD.

Big Data constituie o oportunitate de a utiliza noi tipuri de date în scopul de a crea întreprinderi mai agile, care să rezolve probleme care anterior au fost considerate nesoluționabile, conducând la rezultate mai bune în afaceri. Aceasta duce la schimbări radicale în funcționarea întreprinderilor, care se schimbă de la utilizarea unui model bazat în principal pe experiența decidentului, la un model bazat pe informații, care dă o valoare reală a afacerii și organizației în sine.

#### BIBLIOGRAFIE

1. **DUTCHER, J.:** What Is BD?. Berkeley School of Informatics, Sept. 2014.
2. **GANTZ, J.; REINSEL, D.:** Extracting value from chaos. IDC iView, 2011, pp 1-12.
3. **GARTNER:** "IT glossary: big data" [webpage on the Internet], Stamford, CT; 2012, Available from: <http://www.gartner.com/it-glossary/big-data>.
4. **GARTNER:** 10 Big Data Software Requirements, <http://www.information-management.com/gallery/Big-Data-Required-Software-Applications-10026664-1.html>, accesat august 2015.
5. **HADOOP, A.:** Hadoop, 2009, <http://hadoop.apache.org/>.
6. **HARVEY, C.:** Hadoop and Big Data: 60 Top Open Source Tools, <http://www.datamation.com/applications/hadoop-and-big-data-60-top-open-source-tools-1.html>, iunie 2015.

7. **KHAN, N.; YAQOOB, I.; HASHEM, I. A. T. et al.:** Big Data: Survey, Technologies, Opportunities, and Challenges, The Scientific World Journal, vol. 2014, Article ID 712826, 18 pagini, 2014. doi:10.1155/2014/712826.
8. **ANUGANTI, V.:** Typical “Big” Data Architecture. 2012. Retrieved from: <http://venublog.com/2012/11/30/typical-big-data-architecture/>.
9. **BODAPATI, V.:** Data Integration Ecosystem for Big Data and Analytics. 2013. Retrieved from: <http://smartdatacollective.com/raju-bodapati/103326/data-integration-ecosystem-big-data-and-analytics>.
10. **MARZ, N.; WARREN, J.:** Big data - Principles and best practices of scalable realtime data systems (Chapter 1), 2014.
11. **TIAN, W.; ZHAO, Y.:** Optimized Cloud Resource Management and Scheduling - Theories and Practices, Morgan Kaufmann, Elsevier Inc, 269 p., 2015, ISBN: 978-0-12-801476-9;
12. **MC CREARY, D.; KELLY, D. A.:** Making Sense of NoSQL: A guide for managers and the rest of us, Manning, 2014, ISBN-13: 978-1617291074, ISBN-10: 1617291072.
13. **SADALAGE, P. J.; FOWLER, M.:** NoSQL Distilled: A Brief Guide to the Emerging, World of Polyglot Persistence, [resource.mitfiles.com/CSE/.../NoSQL%20Distilled.pdf](http://resource.mitfiles.com/CSE/.../NoSQL%20Distilled.pdf)