

# APACHE HADOOP

MADIUDIN Radu

Universitatea Tehnică a Moldovei

**Abstract:** *Articolul dat relevă o scurtă introducere în frameworkul Apache Hadoop - framework predestinat lucrului cu masive de informație, oferind o performanță și flexibilitate deosebită. Este descrisă premisa apariției acestuia. Principiile ce stau la baza funcționării reprezintă nucleul frameworkului, constituit din Hadoop HDFS (Hadoop Distributed File System) și MapReduce. Arhitectura aleasă a fost creată în jurul algoritmului MapReduce. Este arătat domeniul utilizării Hadoop și de ce este unul dintre soluțiile viitorului de stocare și prelucrare a informației.*

**Cuvinte cheie:** HADOOP, framework, big data, MapReduce, HDFS, distributed, YARN, Spark, Hive.

## 1. Premise

Ne aflăm în era IOT, rata cantității de informație obținute crește la un nivel comparabil celui exponențial. Domeniile din care informația provine sunt rețele sociale, medicina, secvența genomului, meteo, sateliți, comportamentul cumpărătorilor, informația pe vânzări, finanțe, log file etc. Tot mai multe domenii implementează tehnologiile informaționale cu scopul creșterii productivității creând un aflux puternic de date. Toată informația obținută nu are nici o valoare dacă nu este analizată. Prin big data înțelegem o colecție de datasets ce nu pot fi procesate utilizând tehnici tradiționale de calculare. Aceasta include un volum și o varietate mare ce poate fi împărțită în trei tipuri: informație structurată (informație relațională), informație semi structurată (informație XML), informație nestructurată (Word, PDF, Text, Media Logs).

Metodele tradiționale de soluționare (sistemul centralizat de stocare) ale acestei probleme au devenit inefective. Un sistem centralizat de prelucrare și stocare a informației este efectiv în cazul unui volum mai mic de informație, odată cu creșterea volumului de date apare problema capacității de procesare și network throughput.

## 2. Apache Hadoop

Google a soluționat problema prin utilizarea algoritmului numit MapReduce (figura 1). Acest algoritm divizează taskurile în părți mici și le atribuie la mai multe calculatoare conectate într-o rețea, colectând rezultatul sub forma de final result dataset.

Hadoop este un open source framework Apache scris în Java care a fost creat pe baza principiului sus menționat. Fiind proiectat să fie scalabil de la un singur server la mii de mașini, fiecare oferind stocare și calculare locală, acesta este capabil să permită dezvoltarea unei aplicații ce ar putea efectua o statistică completă la cantități mari de informație. Toate modulele în Hadoop sunt proiectate cu presupunerea fundamentală ca eșecurile hardware se întâlnesc des și trebuie să fie tratate automat de framework.

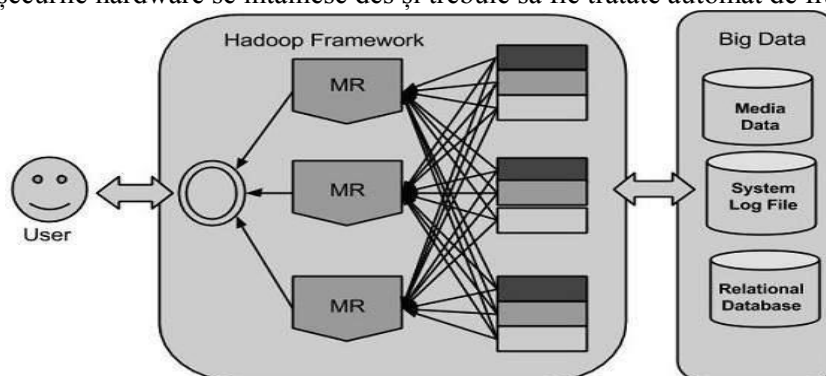


Fig.1. MapReduce

## 3. Hadoop Architecture

Nucleul Apache Hadoop constă din partea storage, cunoscută ca Hadoop Distributed File System (HDFS), și partea procesuală care este modelul programat al algoritmului MapReduce. Frameworkul de baza este compus din următoarele module:

- Hadoop Common: conține librării și utilități necesare funcționării altor module Hadoop. Aceste librării asigură abstracții la nivel filesystem și OS, și conțin java files și scripts necesare pentru startarea Hadoop.
- Hadoop YARN: este un framework pentru job scheduling și cluster resource management.
- Hadoop Distributed File System (HDFS): un sistem distribuit de fișiere care asigură capacitate ridicată de acces la informația aplicației.
- Hadoop MapReduce: este un sistem YARN utilizat în procesarea paralelă a seturilor largi de informație.

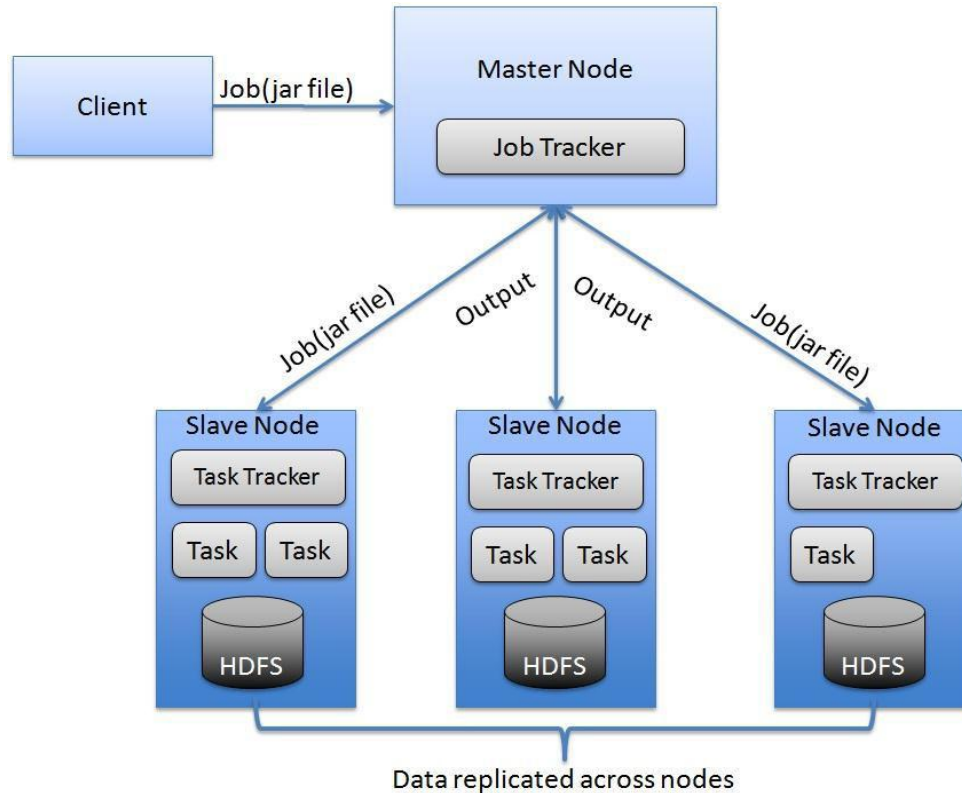


Fig. 2. Procesarea informației

#### 4. Cum funcționează Hadoop?

##### Stage 1:

Un utilizator/aplicație poate înainta un job către Hadoop (figura 2) pentru procesul necesar prin specificarea următoarelor puncte:

- locația fișierelor input și output în sistemul distribuit.
- clasele Java sub formă de fișiere jar ce conțin implementarea funcțiilor map și reduce.
- configurația job prin setarea diferitor parametri specifici jobului.

##### Stage 2:

Clientul job Hadoop înaintază jobul (jar/executabil) și configurația către JobTracker care preia responsabilitatea distribuirii softului/configurației la slaves, planificarea taskurilor și monitorizarea lor, furnizând statusul și diagnosticul informației către job-client.

##### Stage 3:

TaskTrackers pe noduri diferite execută taskurile conform implementării MapReduce și rezultatul funcției reduce este salvat în fișierele output pe sistemul de fișiere (figura 3, 4).

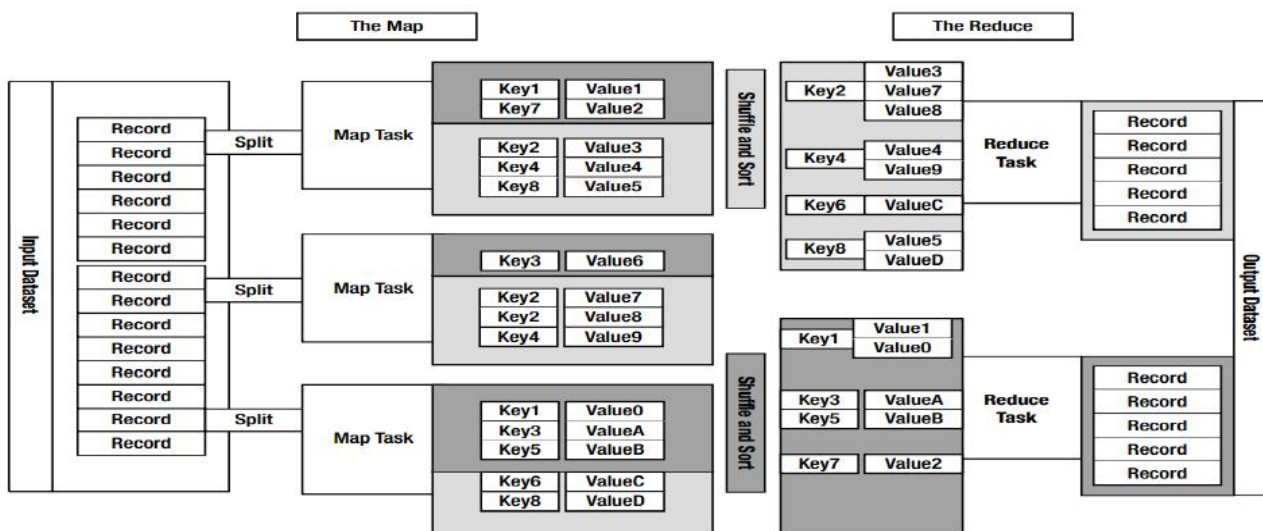


Fig. 3. Principiul de funcționare a frameworkului

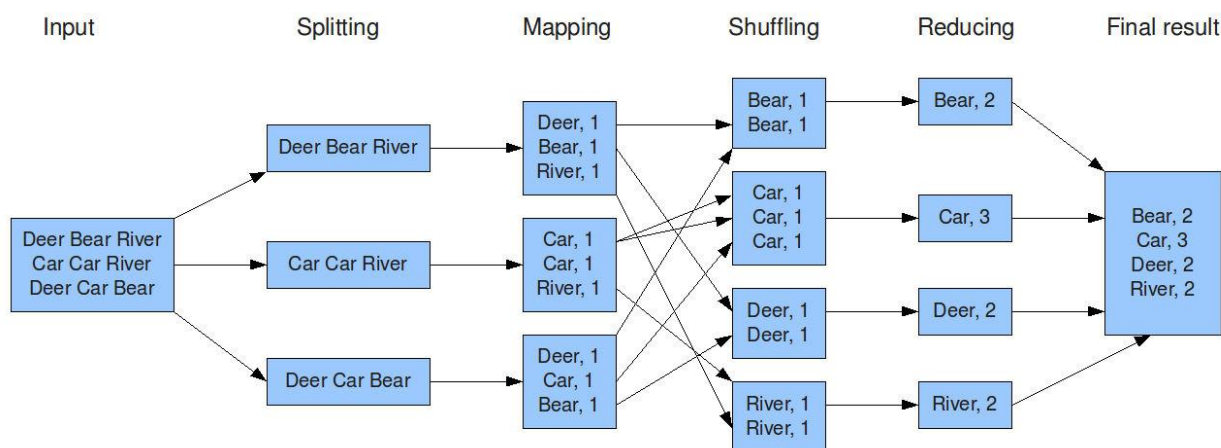


Fig. 4. Exemplu practic a algoritmului MapReduce

## 5. Hadoop Ecosystem

Ecosistemul Hadoop include atât proiecte Apache open source cât și o gama largă de instrumente și soluții comerciale (figura 5). Unele dintre cele mai cunoscute open source sunt Spark, Hive, Pig, Oozie și Sqoop. Opțiunile comerciale sunt încă mai diverse, incluzând platforme și distribuții de la vânzătorii ca Cloudera, MapR, plus o varietate de instrumente specifice sarcinilor Hadoop de dezvoltare, producere și mentenanță

Spark este atât un model de programare cât și unul computațional. Acesta oferă o poartă la computarea în memorie pentru Hadoop, ce reprezintă unul din motivele popularității lui și o adopție largă. Spark oferă o alternativă la MapReduce ce permite executarea sarcinilor în memorie, în schimb pe disk. Utilizând computarea în memorie, volumul de lucru rulează de la 10 la 100 de ori mai repede în comparație cu execuția pe disk.

Hive reprezintă un soft de depozitare a datelor care răspunde de faptul cum informația este structurată și interogată în clusterelor distribuite Hadoop. Hive este, de asemenea, un mediu popular utilizat în scrierea interogărilor pentru informația în mediul Hadoop. Oferă instrumente ETL (Extract, Transform and Load) ce aduc mediului capacități asemănătoare SQL.

Pig este un limbaj procedural utilizat în dezvoltarea aplicațiilor procesate în paralel pentru seturi mari de data în mediul Hadoop. Pig este o alternativă programării Java. Pig include Pig Latin, ce este un scripting limbaj. Pig traduce scripturile Pig Latin în MapReduce, ce apoi pot fi rulate pe YARN și procesarea informației în clusterelor HDFS. Este popular din motivul automatizării unei părți complexe a dezvoltării MapReduce.

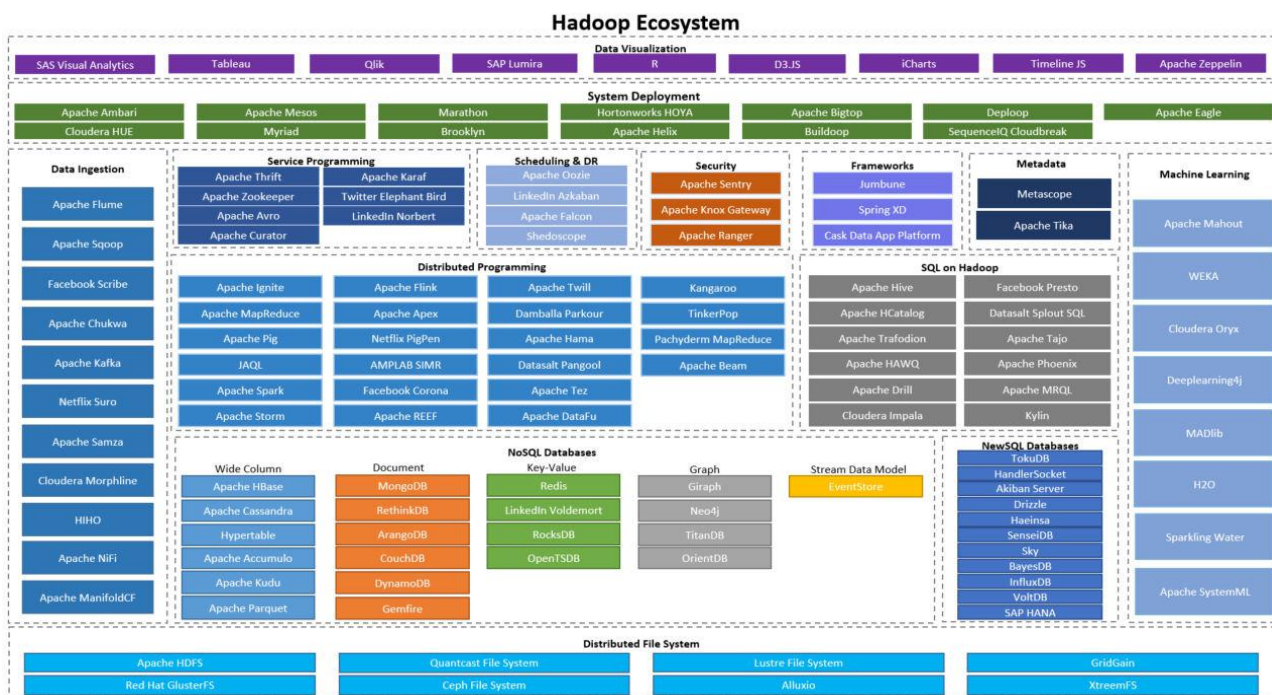


Fig. 5. Familia Hadoop

## 6. Concluzii

Hadoop este una dintre cele mai performante și larg utilizate sisteme în manipularea volumelor mari de informație rapid și ieftin. Acesta permite o flexibilitate deosebită prin diversitatea modulelor oferite, o posibilă scalabilitate lineară fără pierderea vitezei de analiză și procesare. Oferă o redundanță a datelor datorită stocării informației în clustere de mașini. Datorită acestor calități Hadoop a fost pe larg implementat de marii giganti software ca Google, Amazon, Facebook etc. Acest framework a devenit o parte a dezvoltării curentului IOT facilitând și implementând aspectul descentralizat al acestuia, oferind suportul necesar funcționării eficiente.

## Bibliografie

1. Alex Homes. *Hadoop în Practice*. Second Edition. Publisher: Manning Publications, September 2014, ISBN: 9781617292224.
2. Jonathan R. Owens, Jon Lentz, Brin Femiano. *Hadoop Real World Solutions Cookbook*. [Resursa electronica] – Regim de acces: <https://www.amazon.com/Hadoop-Real-World-Solutions-Cookbook-ebook/dp/B00AIVQE3I>
3. Chuck Lam. *Hadoop în Action*. Publisher: Manning Publications. December 2010, ISBN: 9781935182191.
4. *Tutorialspoint* [Resursa electronica] – Regim de acces: <https://www.tutorialspoint.com/hadoop>
5. Apache\_Hadoop#Architecture, Wikipedia [Resursa electronica] – Regim de acces: [https://en.wikipedia.org/wiki/Apache\\_Hadoop#Architecture](https://en.wikipedia.org/wiki/Apache_Hadoop#Architecture)
6. Hadoop Ecosystem and Components [Resursa electronica] – Regim de acces: <http://www.bmc.com/guides/hadoop-ecosystem.html>