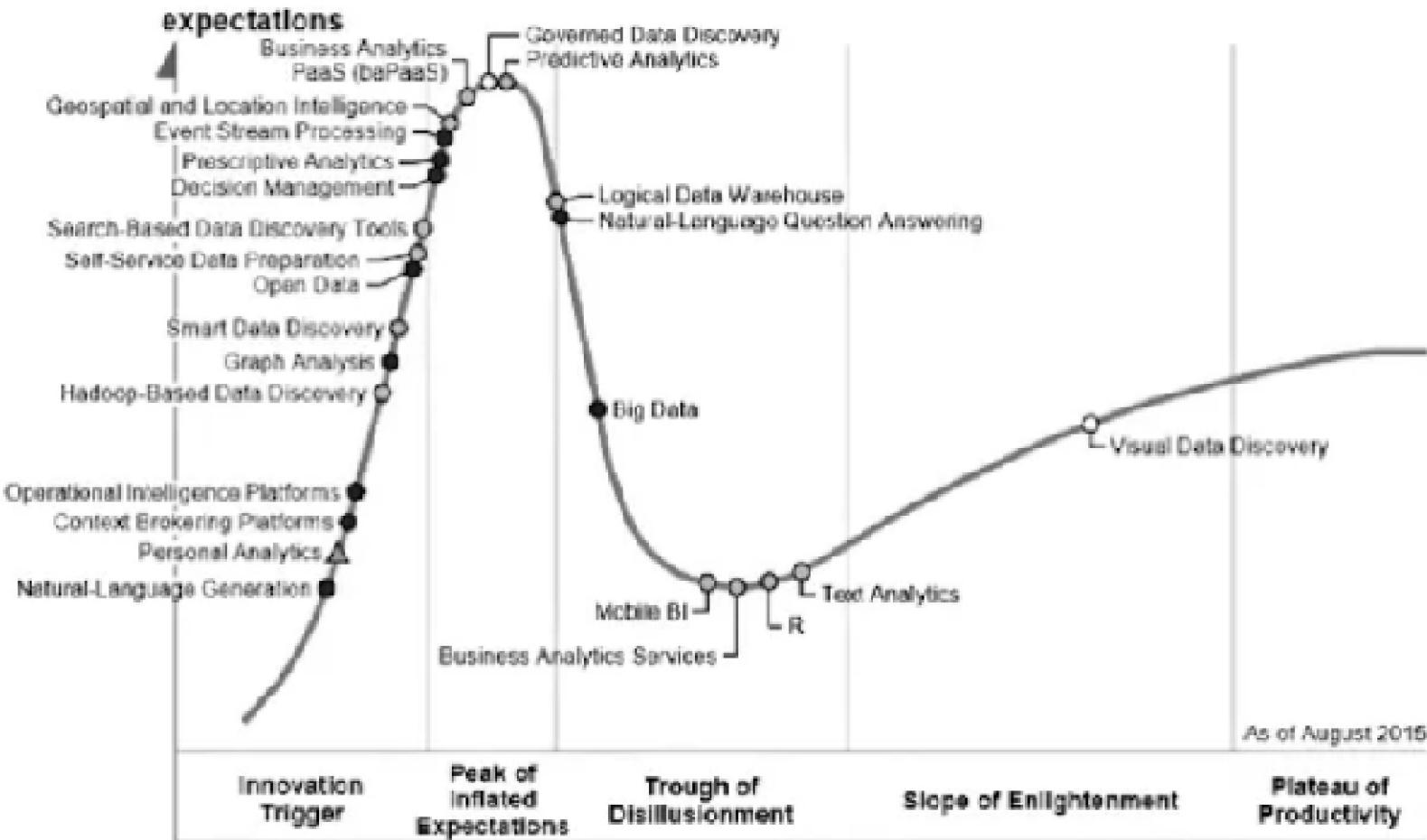


# **Sisteme Distribuite**

Cursul 12

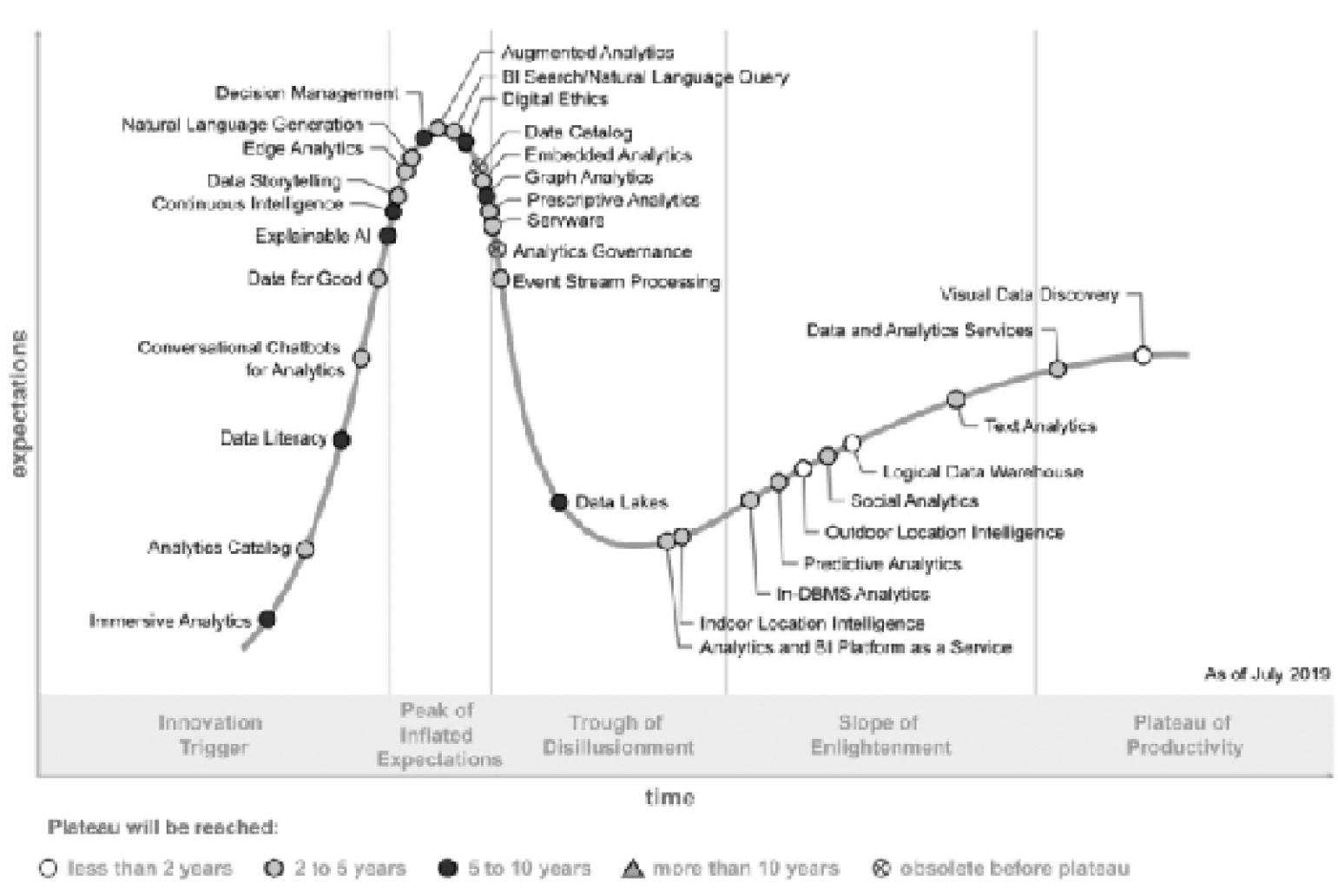
Mihai Zaharia

# De ce?



Gartner - hyper cicle for bussiness intelligence and analitics 2015

# De ce?



Gartner - hyper cicle for bussiness intelligence and analitics 2019

# **Bussiness Intelligence**

- X

# **Bussiness Analytics**

user=Corporate Manager

# **Datele în business analytics&intelligence**

- X

# **Data governance -> bigdata**

- X

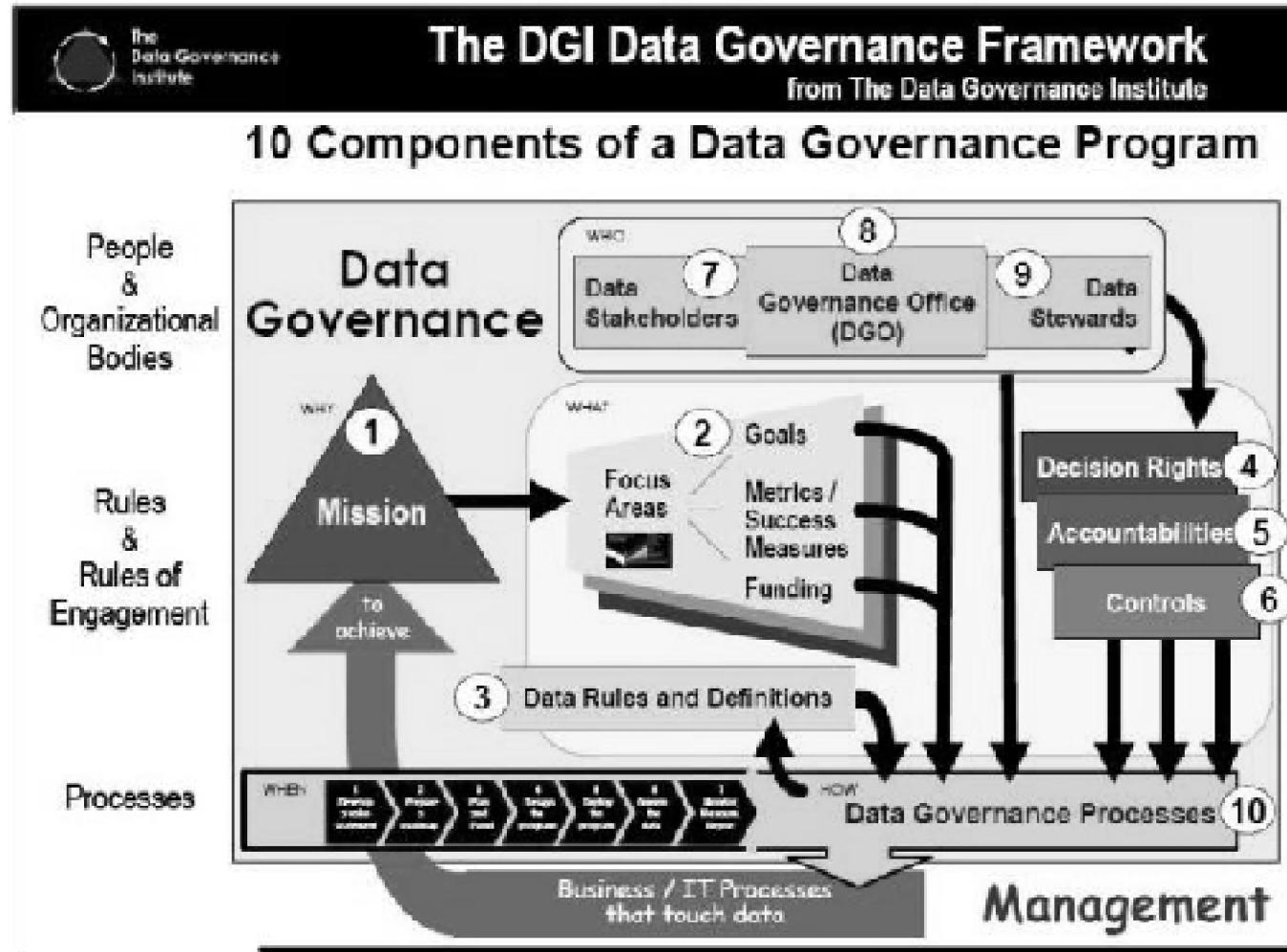
# **Guvernarea datelor - istoric**

- X

# **Guvernarea datelor**

- X

# Guvernarea datelor



- X

# **Importanța guvernării datelor**

- X

# **Scopurile guvernării datelor**

- X

# **GDPR vs guvernarea datelor**

- X

# **Responsabilități pentru guvernarea datelor**

cauza morții CDO-ului!

# **Responsabilități pentru guvernarea datelor**

Managerul și echipa pentru guvernarea datelor

# **Responsabilități pentru guvernarea datelor**

Comitetul

# **Responsabilități pentru guvernarea datelor**

Data stewards

# **Guvernarea big data**

depinde de nivelul din care este analizată

# **Calitatea datelor**

cum se face?!

# **Gestiunea metadatelor**

- X

**Chiar este înțeleasă guvernarea datelor?**

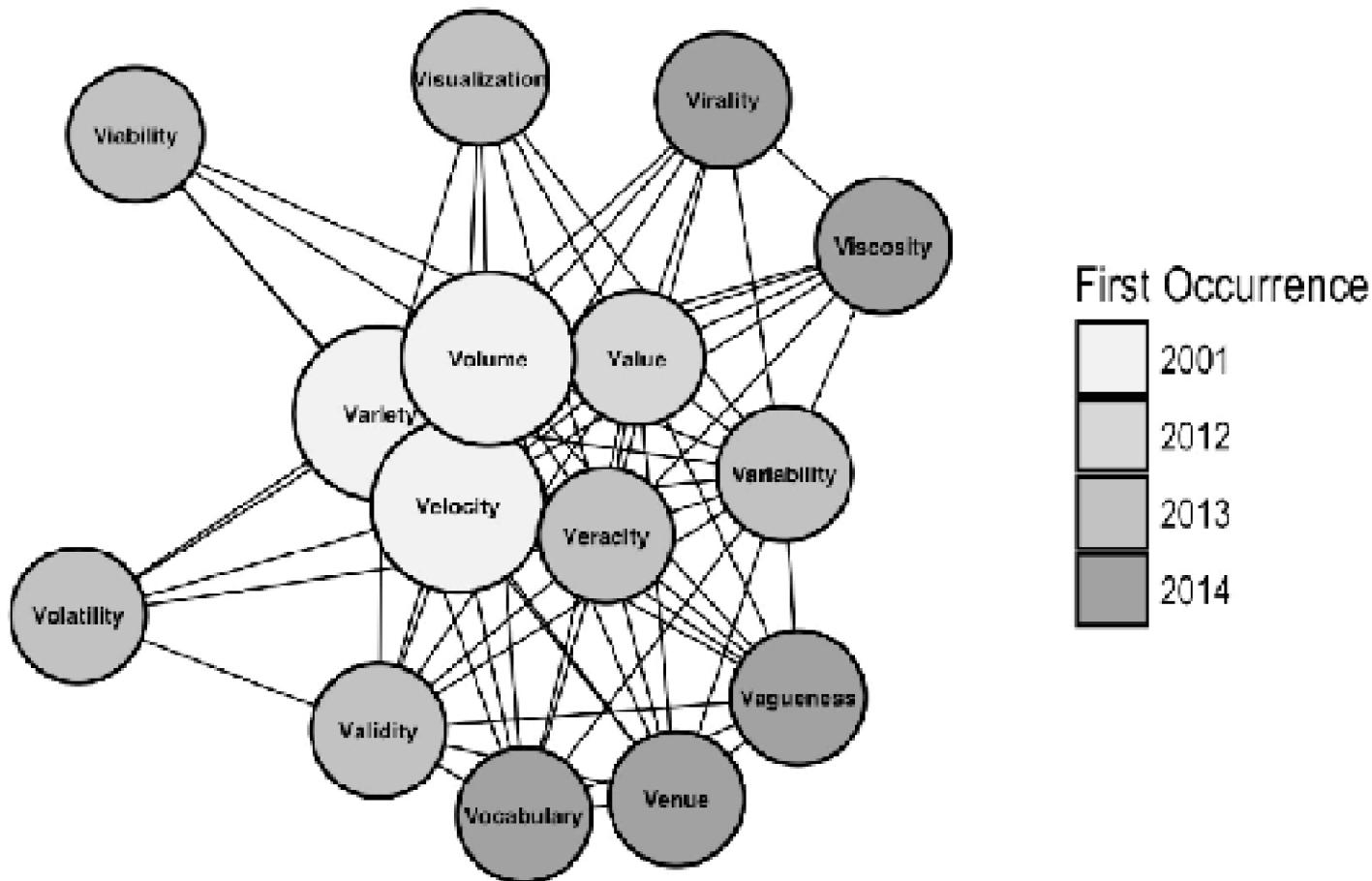
Oh my God It's full of stars!

# Big Data

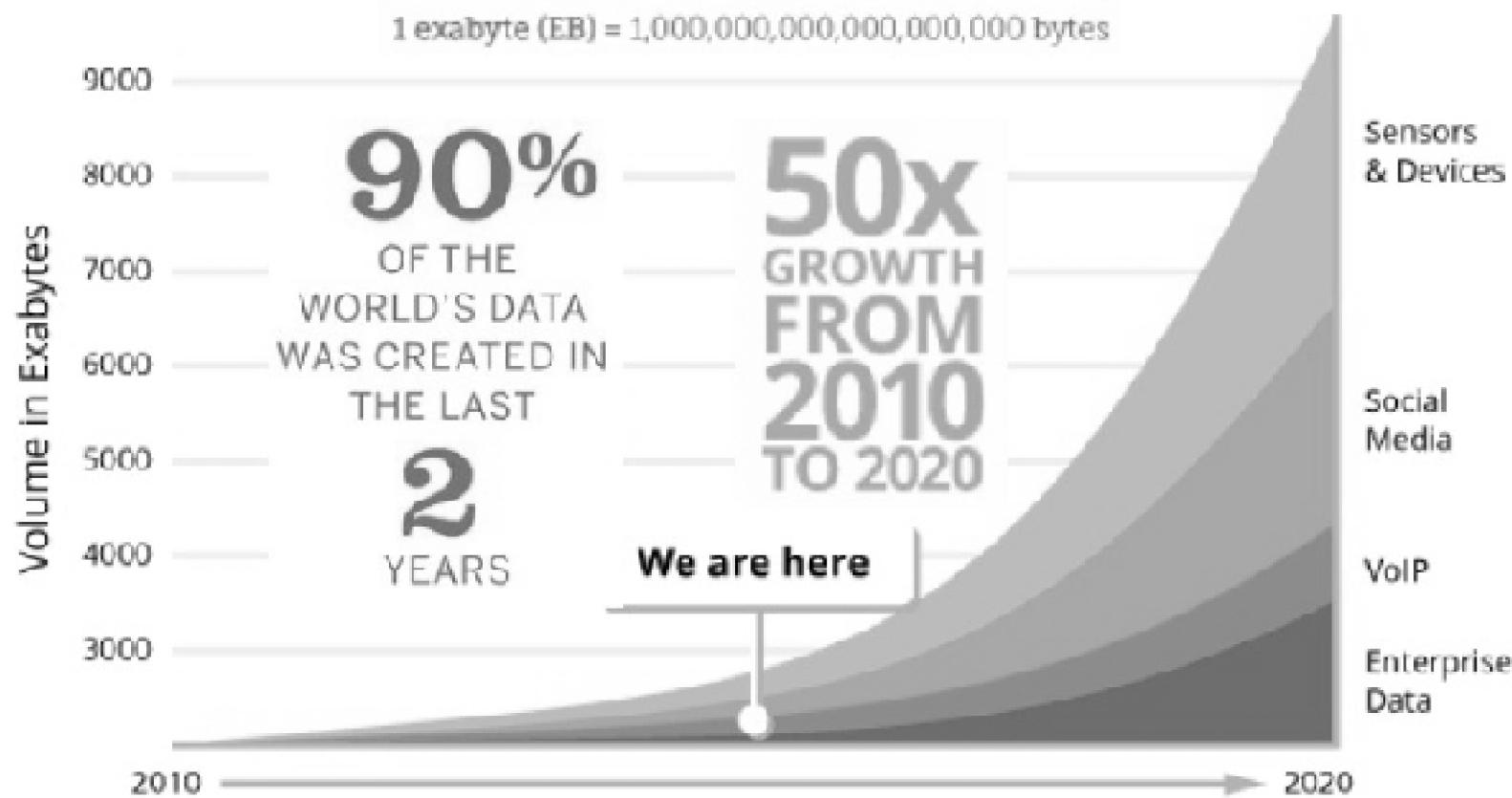
- *"Visualization provides an interesting challenge for computer systems: data sets are generally quite large, taxing the capacities of main memory, local disk, and even remote disk. We call this the problem of big data. When data sets do not fit in main memory (in core), or when they do not fit even on local disk, the most common solution is to acquire more resources."*

Cox si Ellsworth, IEEE, 1997

# Dimensiunile big data

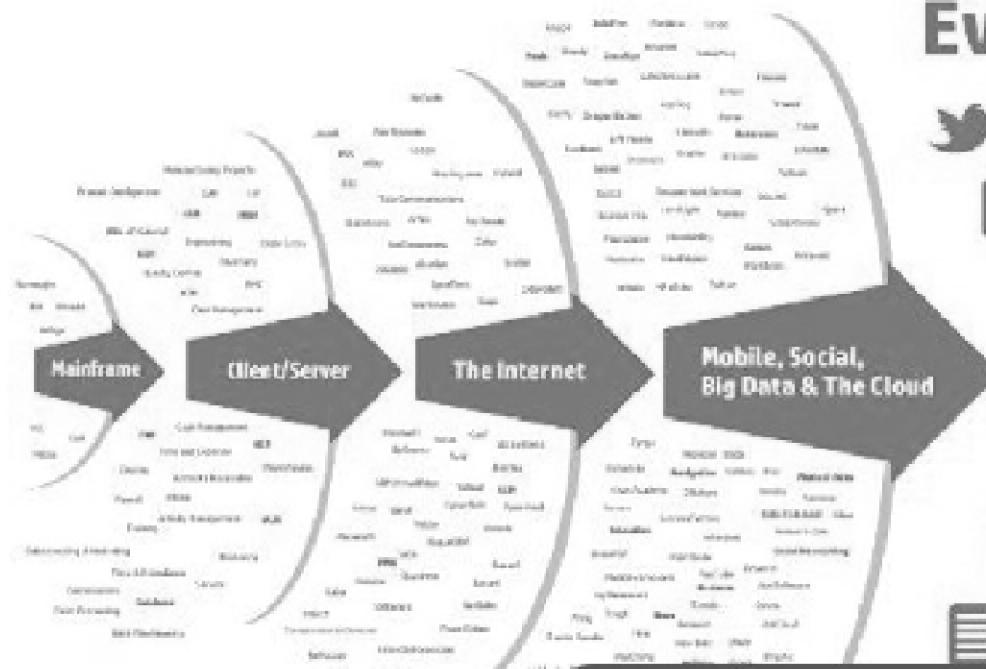


# Dimensiunile Big Data 1. Volumul



creștere exponențială a datelor generate

# Dimensiunile Big Data 2. Velocitate



**Every 60 seconds**

**98,000+ tweets**

**695,000 status updates**

**11million instant messages**

**698,445 Google searches**

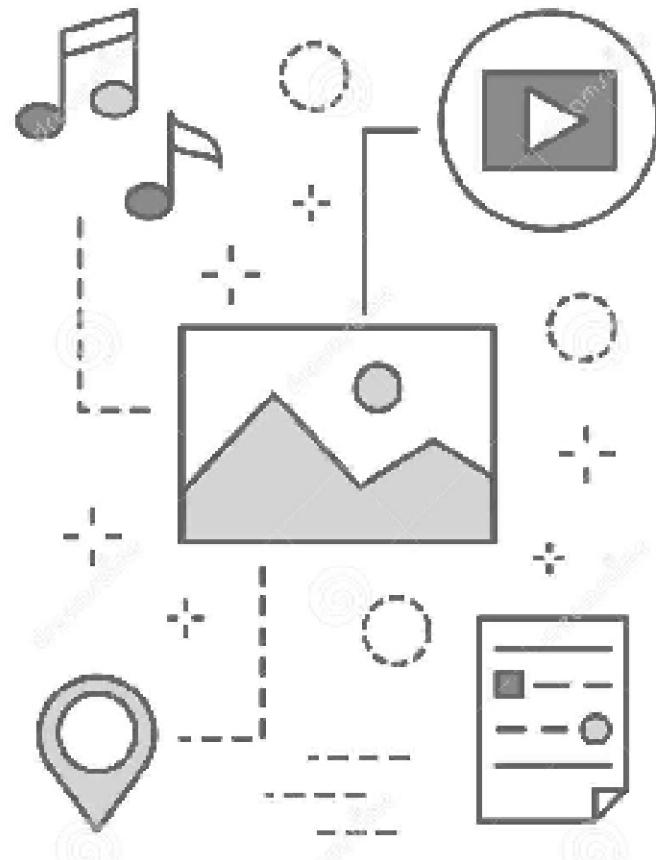
**168 million+ emails sent**

**1,820TB of data created**

**217 new mobile web users**

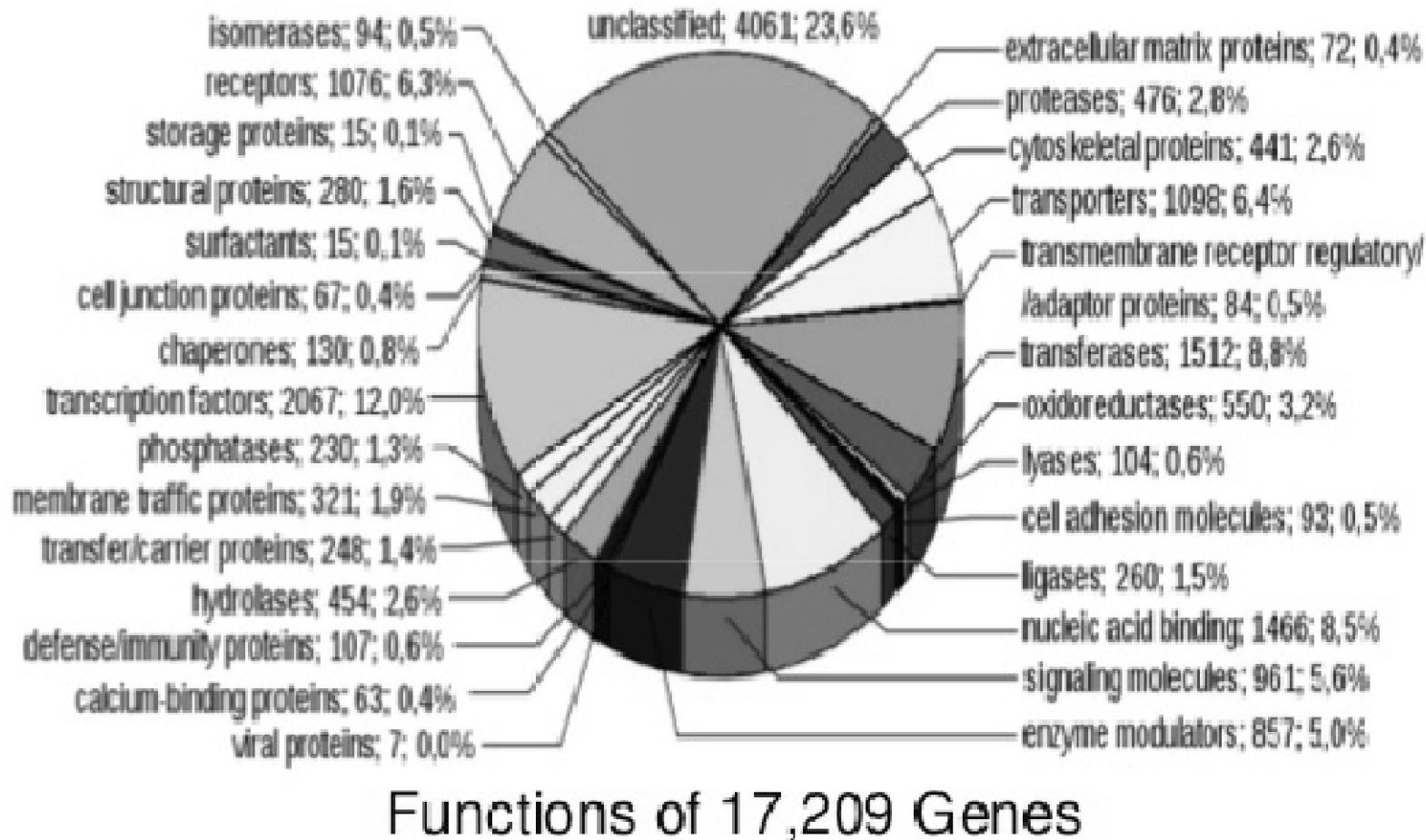
viteza de producție a noilor date

# Dimensiunile Big Data 3. Varietate



- Datele structurate
- Datele semistructurate
- Date nestructurate

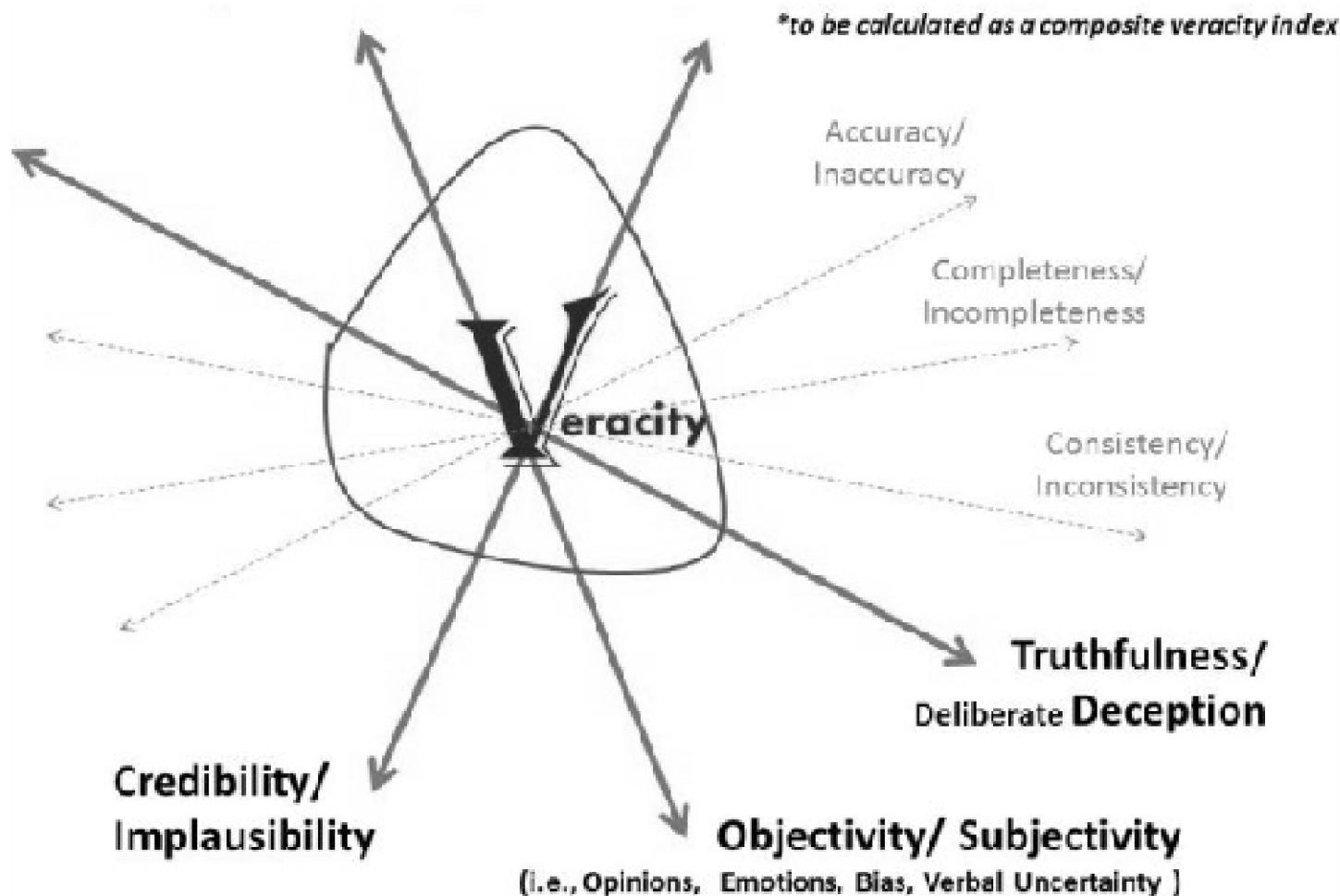
# Dimensiunile Big Data 4. Variabilitate



Functions of 17,209 Genes

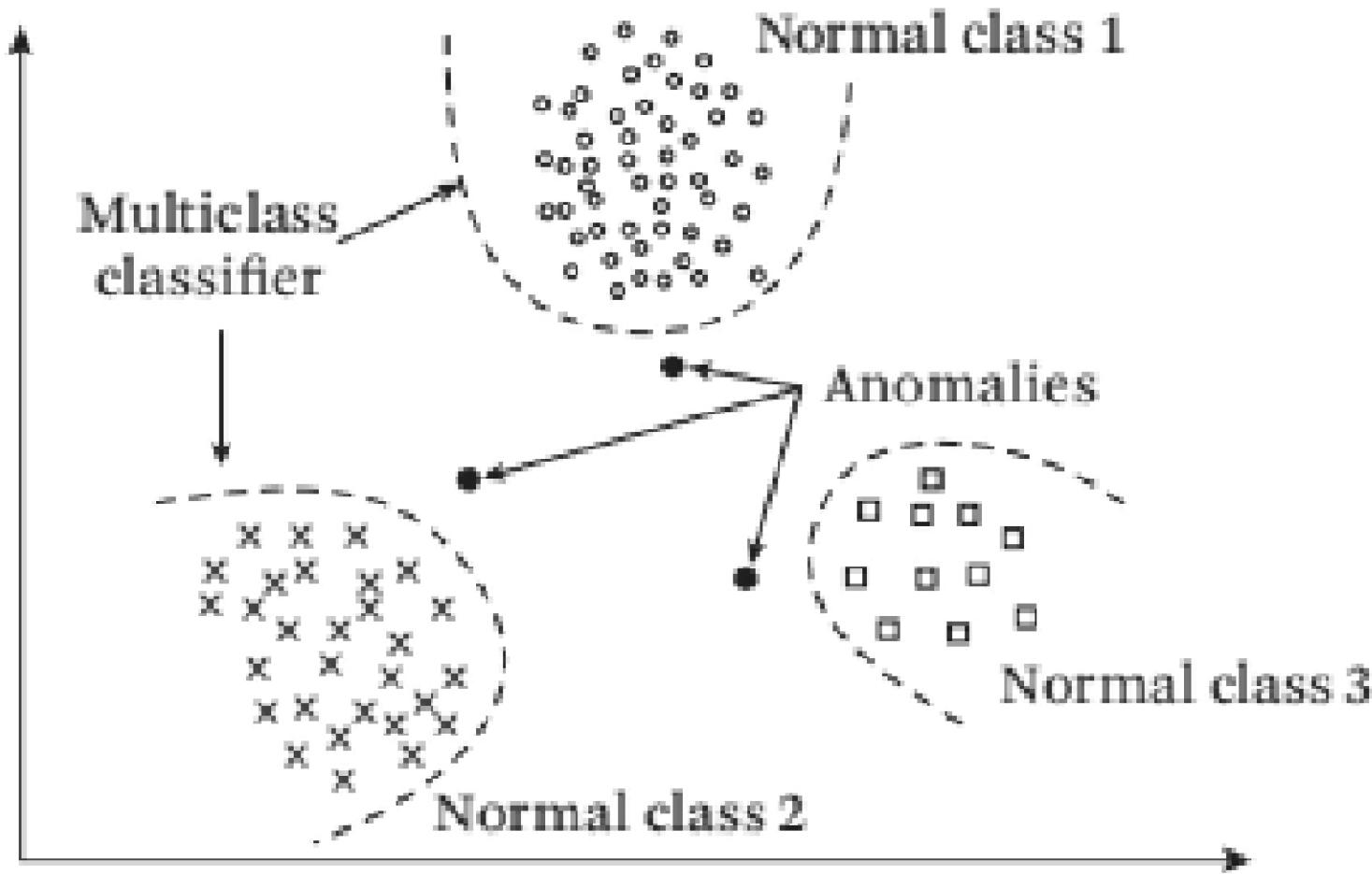
Neomogenitatea sau  
Variatia dimensională și compozițională

# Dimensiunile Big Data 5. Veracitate



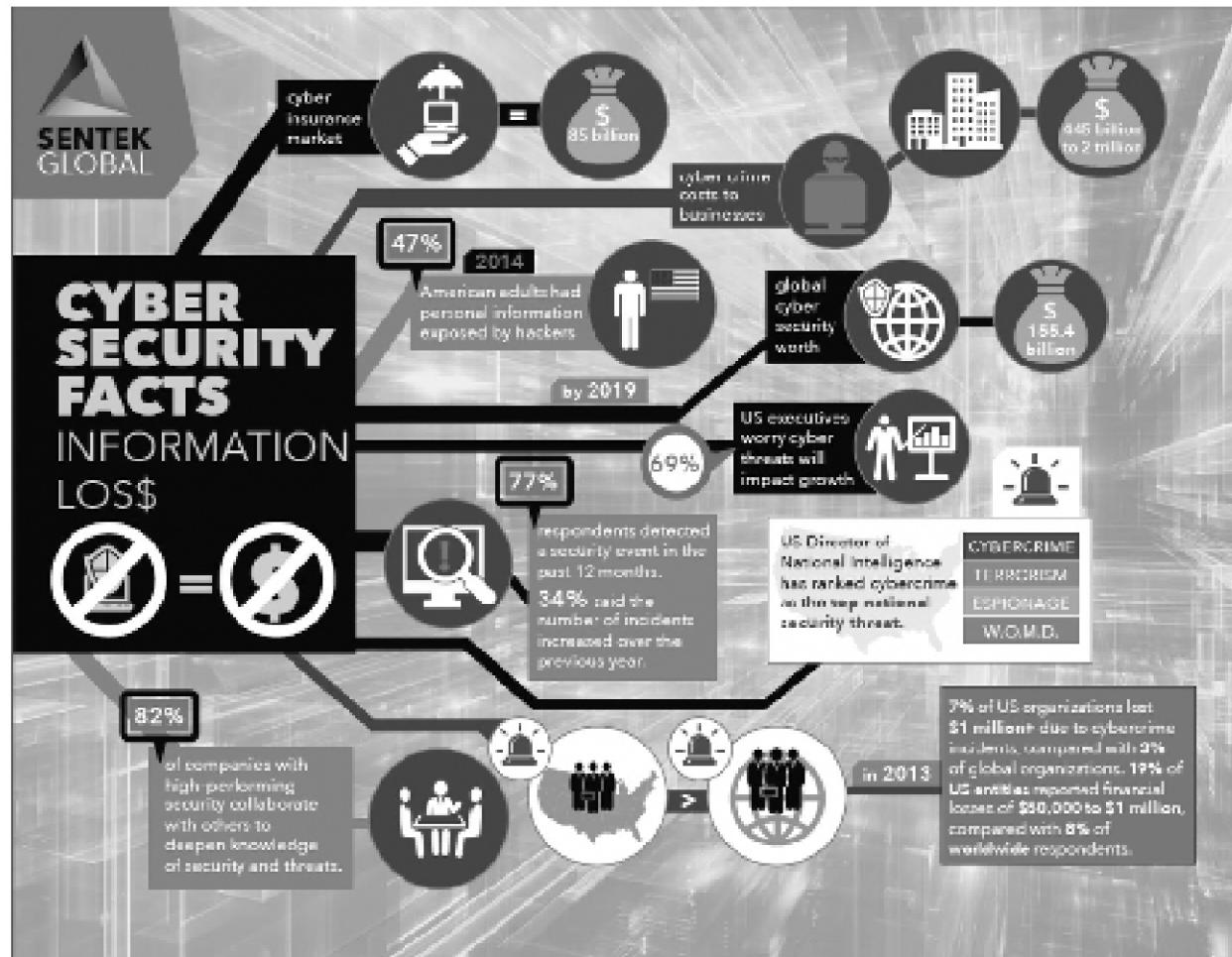
veriditate sau corectitudine

# Dimensiunile Big Data 6. Validitatea



Detectarea anomaliilor după gruparea pe categorii

# Dimensiunile Big Data 7. Vulnerabilitate



după Forbes

# **Dimensiunile Big Data 8. Volatilitate**

date noi vs. date vechi?

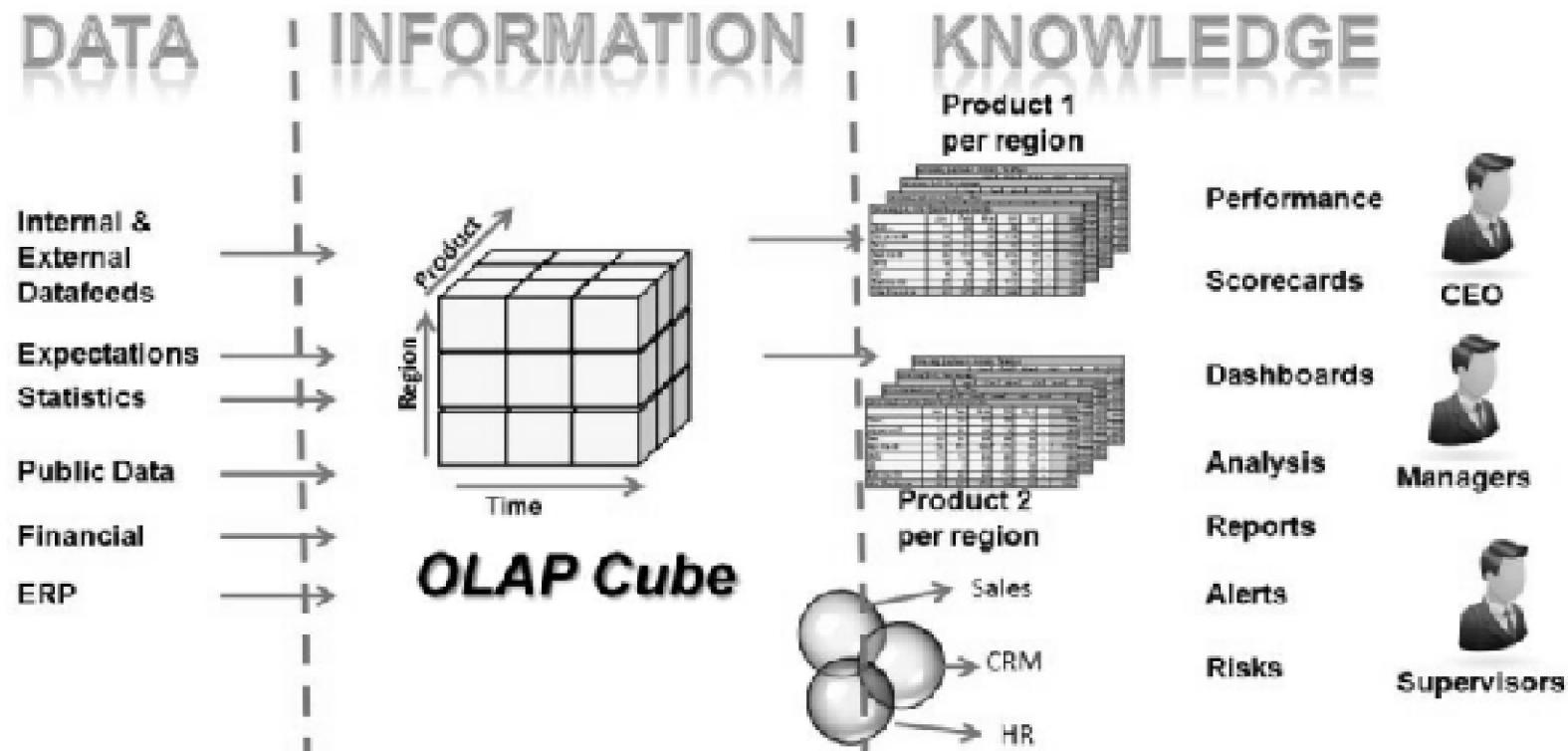
# **Dimensiunile Big Data 9. Vizualizarea**

Google Facets

# **Dimensiunile Big Data 10. Valoare**

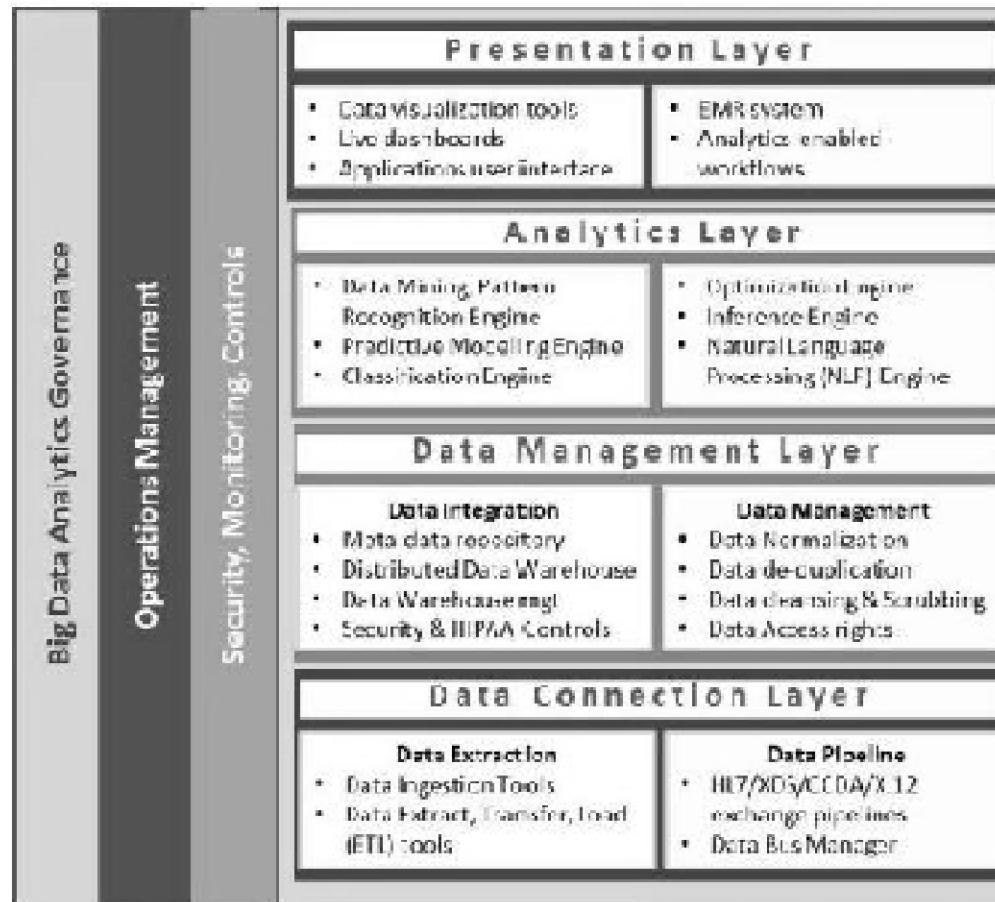
guvern, corporații, IMM, PFA

# Cubul OLAP



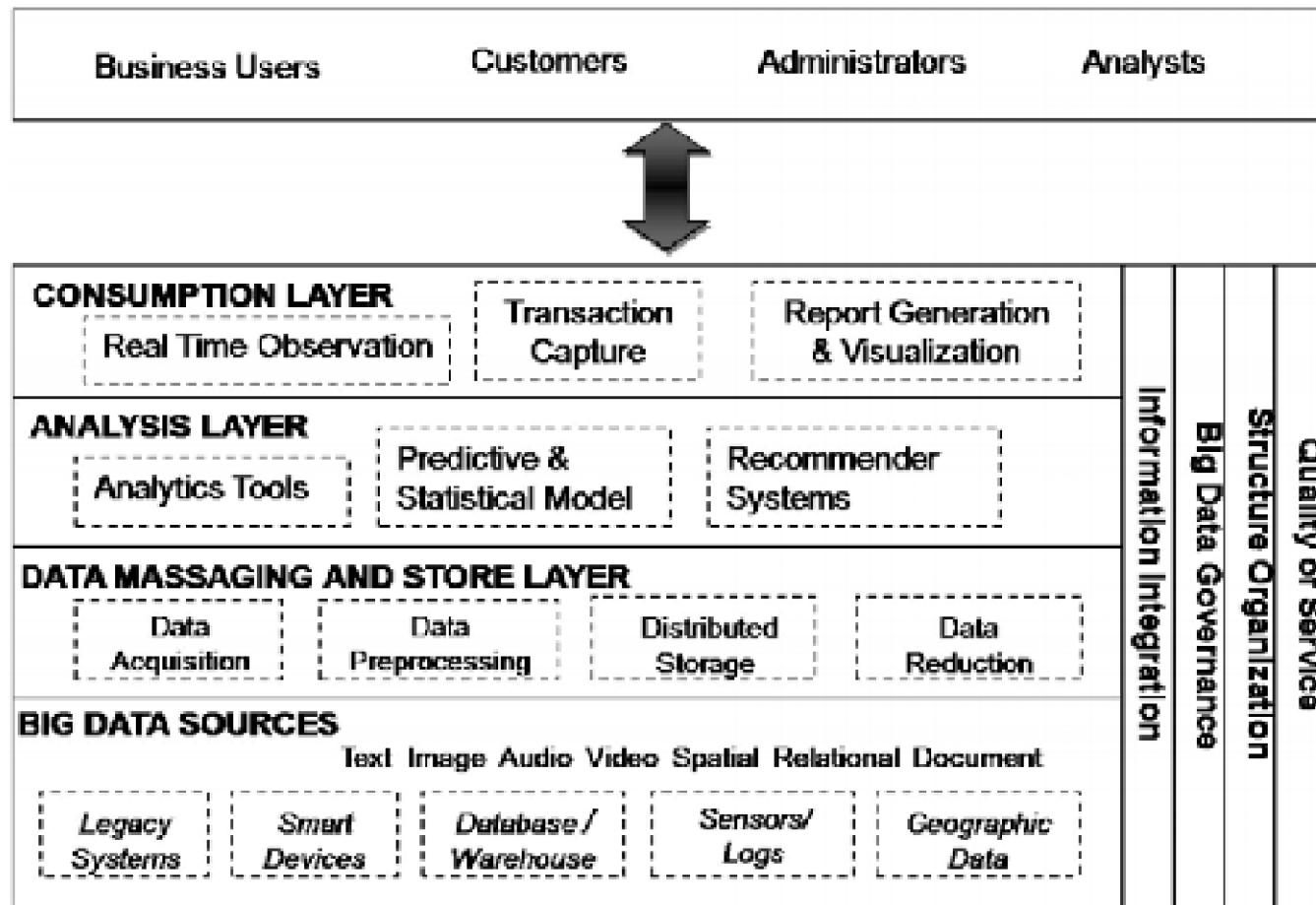
- d

# Arhitectura unei platforme



pentru data analitics

# Arhitectura unei platforme



- X

# Big data governance

Organization	Establish Data Council, Data Stewards, Data Governance Steering Committee
Metadata	Data definitions, data lineage, technical metadata, data registration
Compliance	Regulatory audits, policies, compliance with security/privacy policies
Data Quality	Ensure data is complete & correct. Measure, improve, certify data.
Business Process Integration	Data procurement, ownership, policies around data frequency, availability, etc.
Master Data Management	Establish data & usage taxonomy: business critical hierarchy, Admirals, Members, Providers, Consumers, etc.
Information Lifecycle Management (ILM)	Data retention, regulatory compliance with data/model snapshots; purge Schedule, storage/archiving

componente

# Big data governance

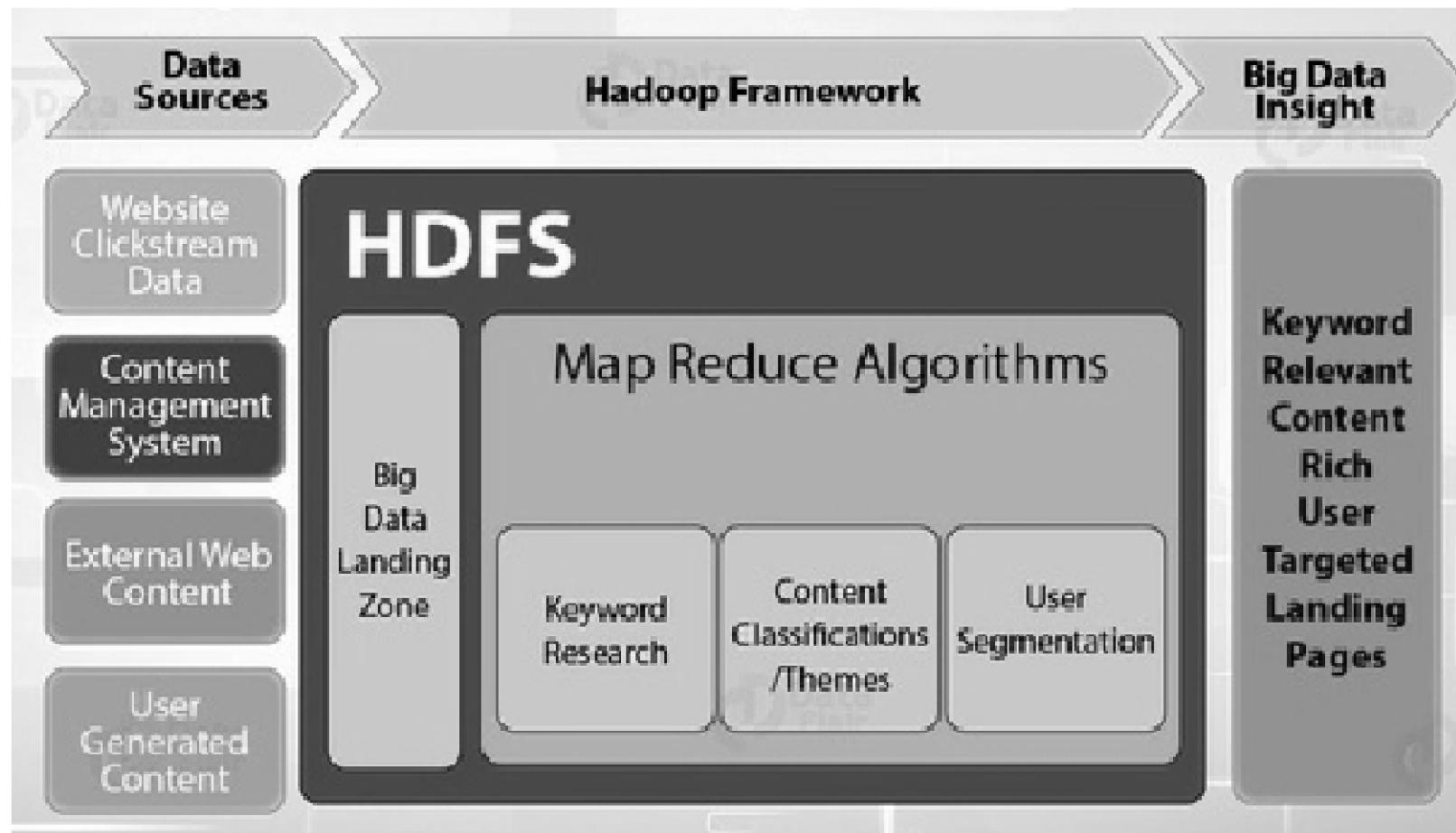


bazată pe Hadoop

# **Hadoop?**

studenți sau profesori?

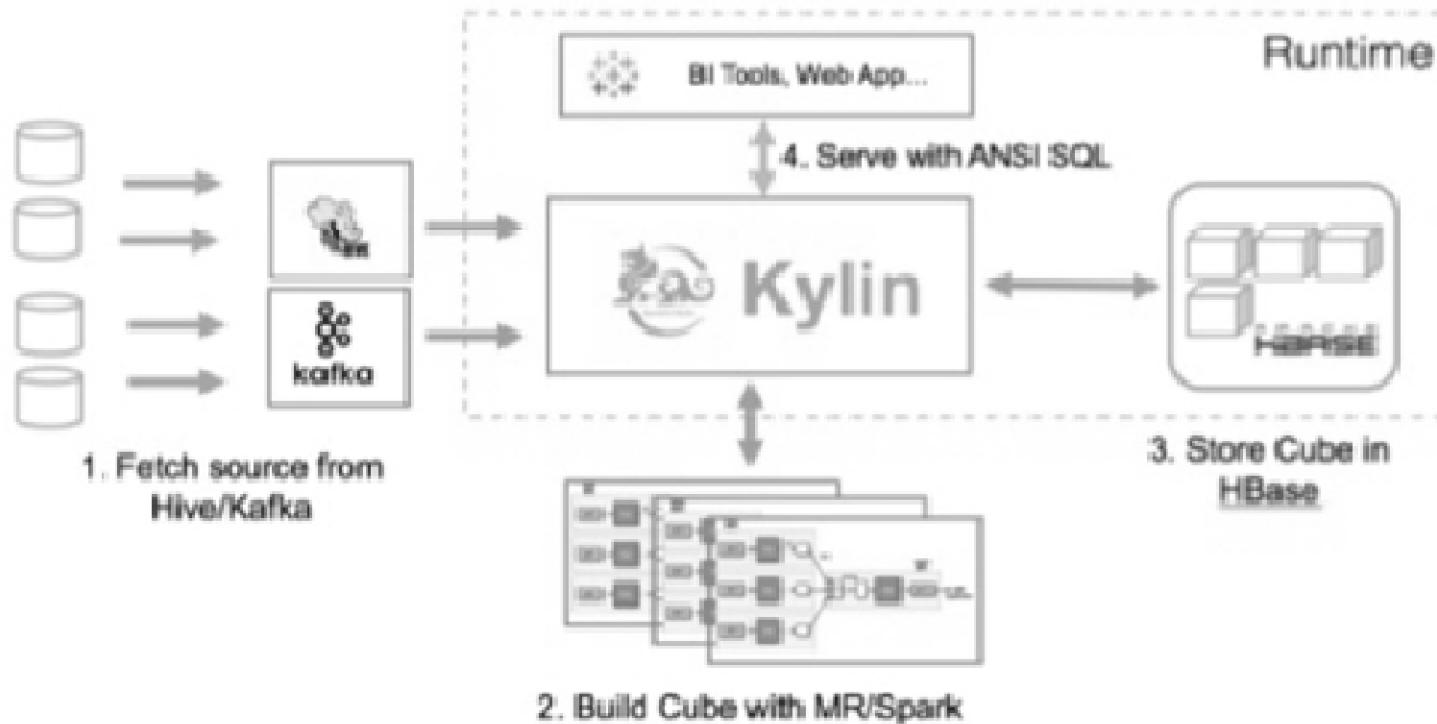
# Hadoop



- X

# **Guvernarea datelor: Tactică sau Strategică?**

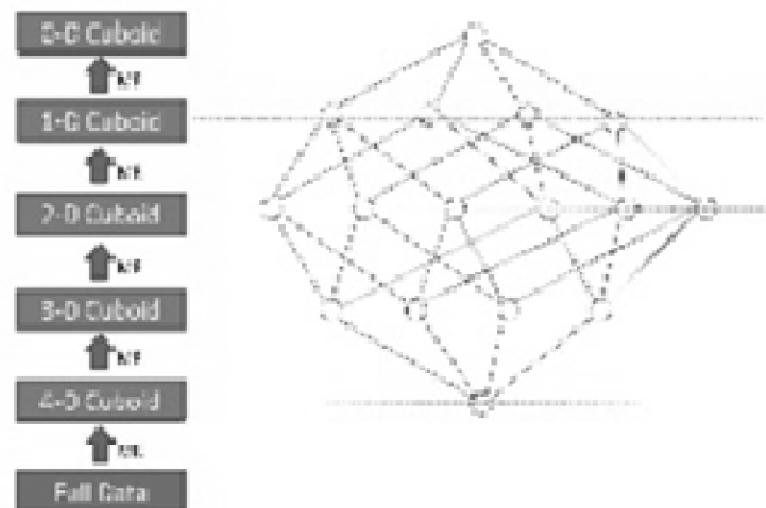
# OLAP & Hadoop



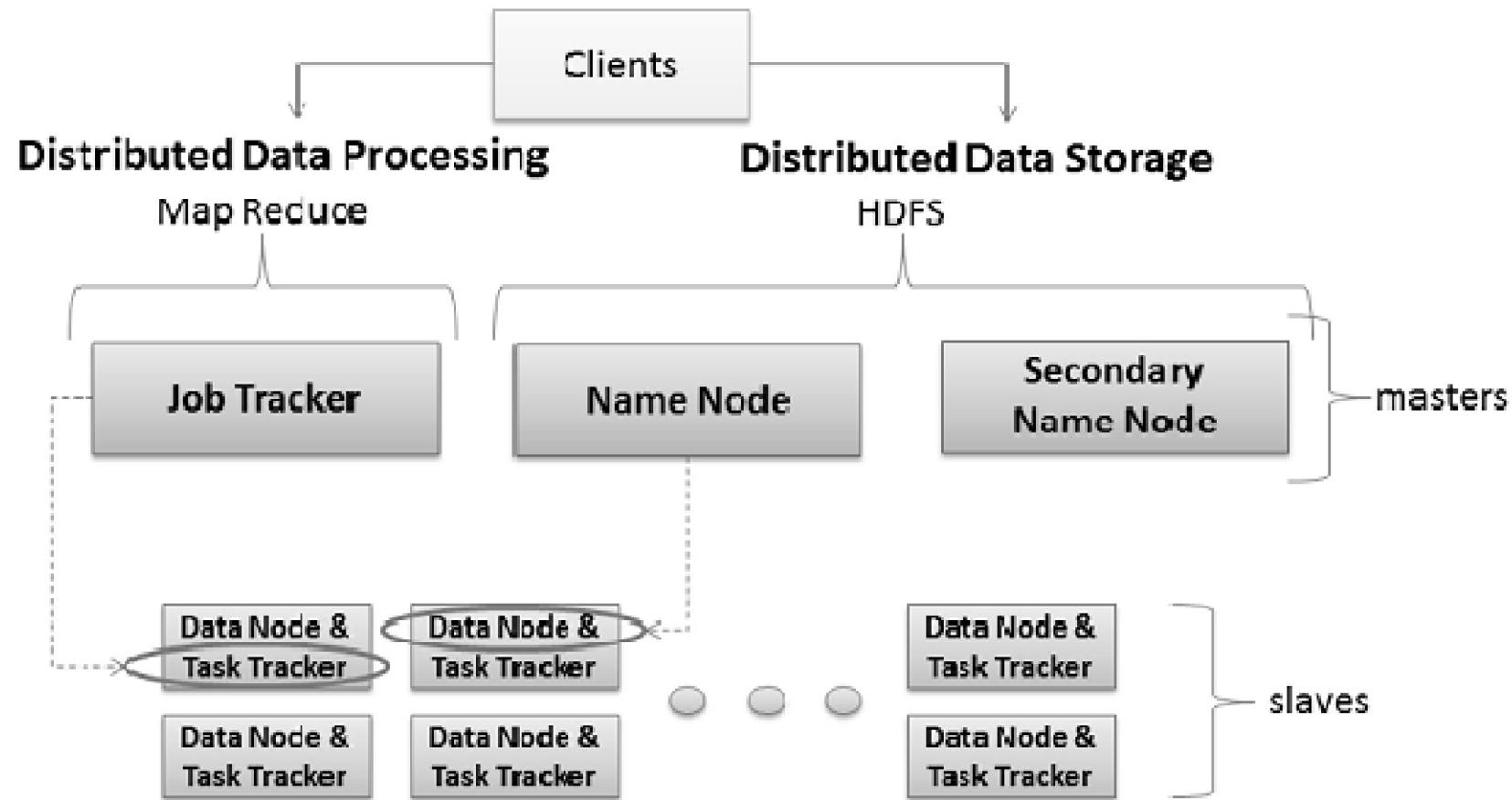
- un exemplu bazat pe kylin

# OLAP & Hadoop

- Calculate Cuboids by layer : N dim (Base cuboid), N-1 dim, N-2..., 1, 0
- Reuse previous layer's result
- HDFS used for data sharing
- Totally need N round MR;
- generarea cubului cu mapreduce

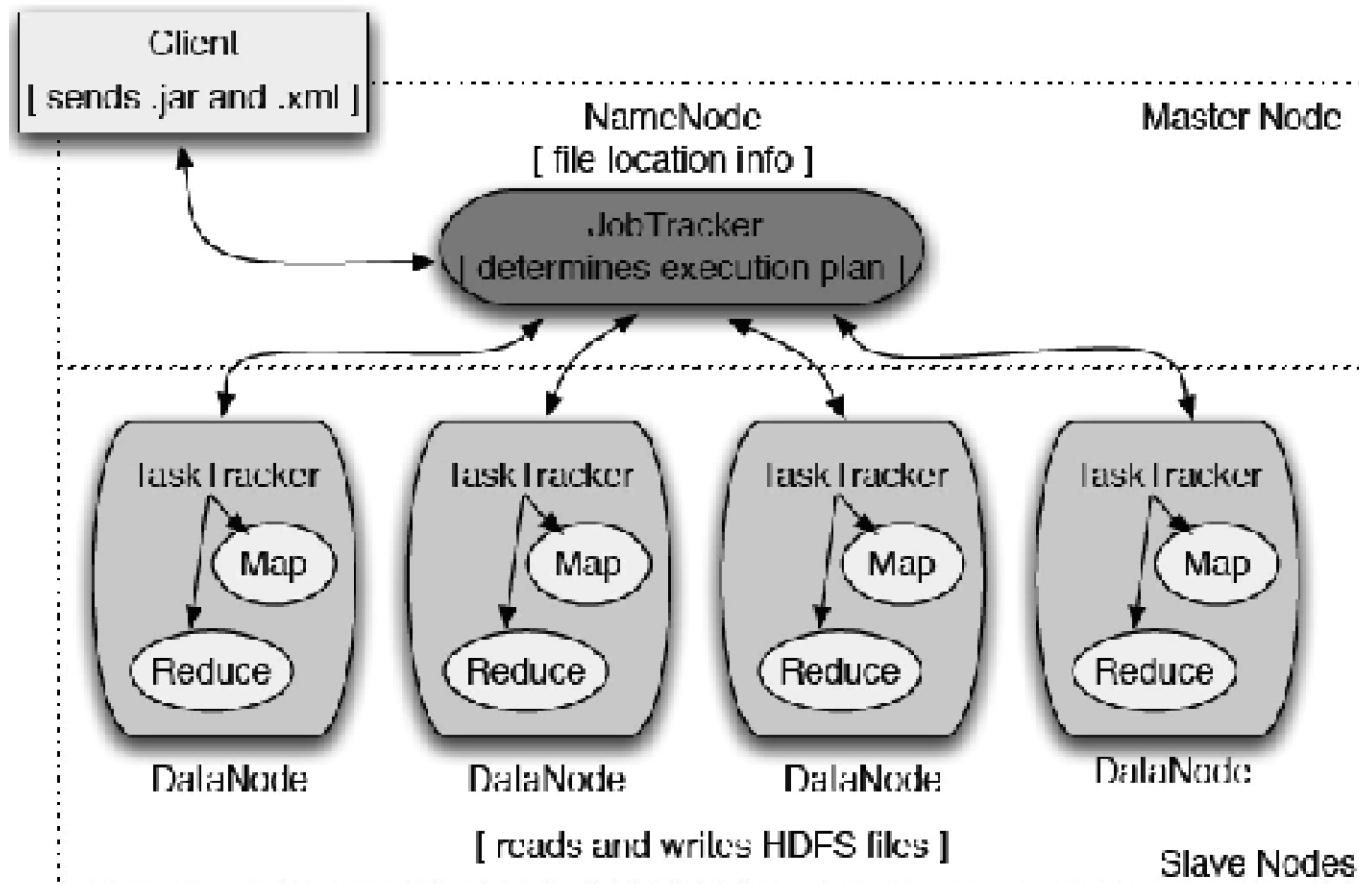


# HDFS



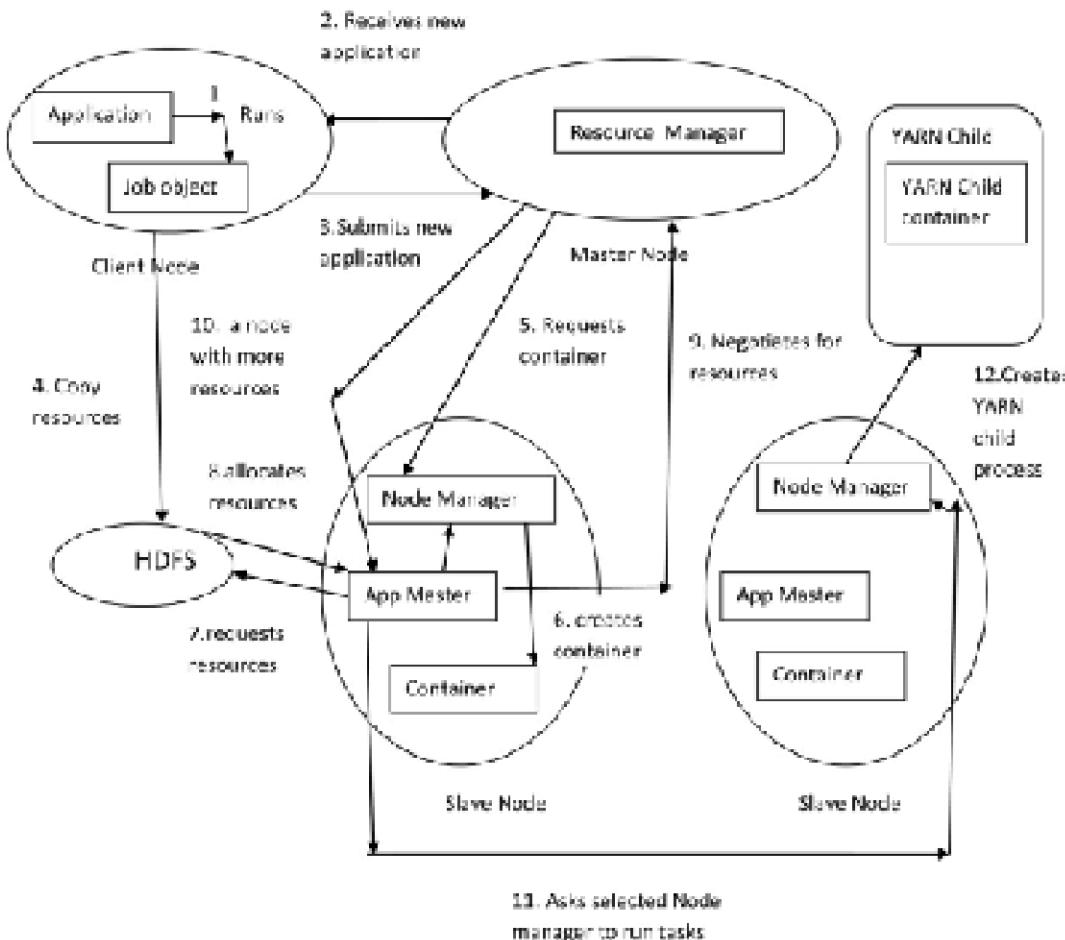
- roluri server in hadoop

# MapReduce



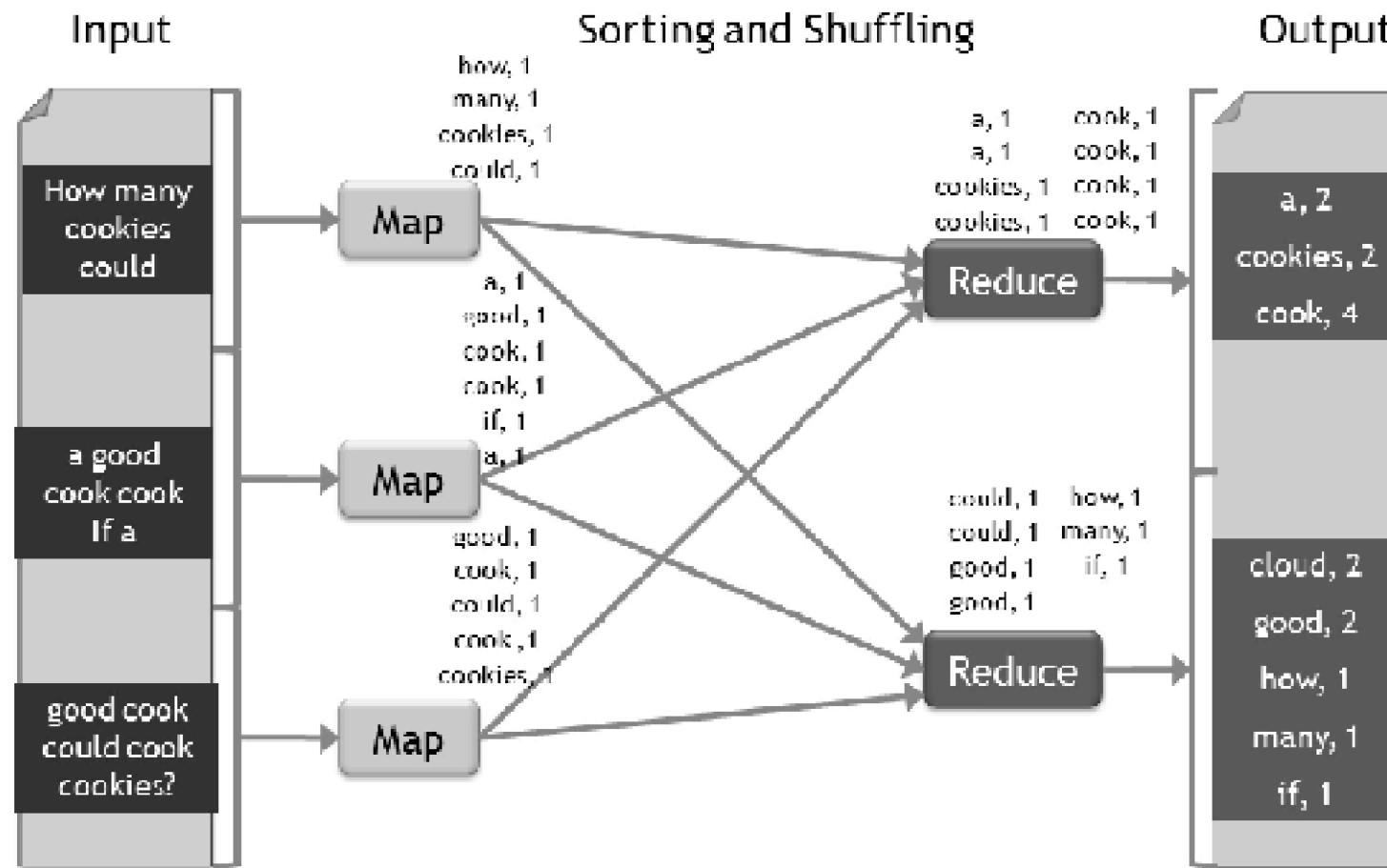
un instrument!

# MapReduce



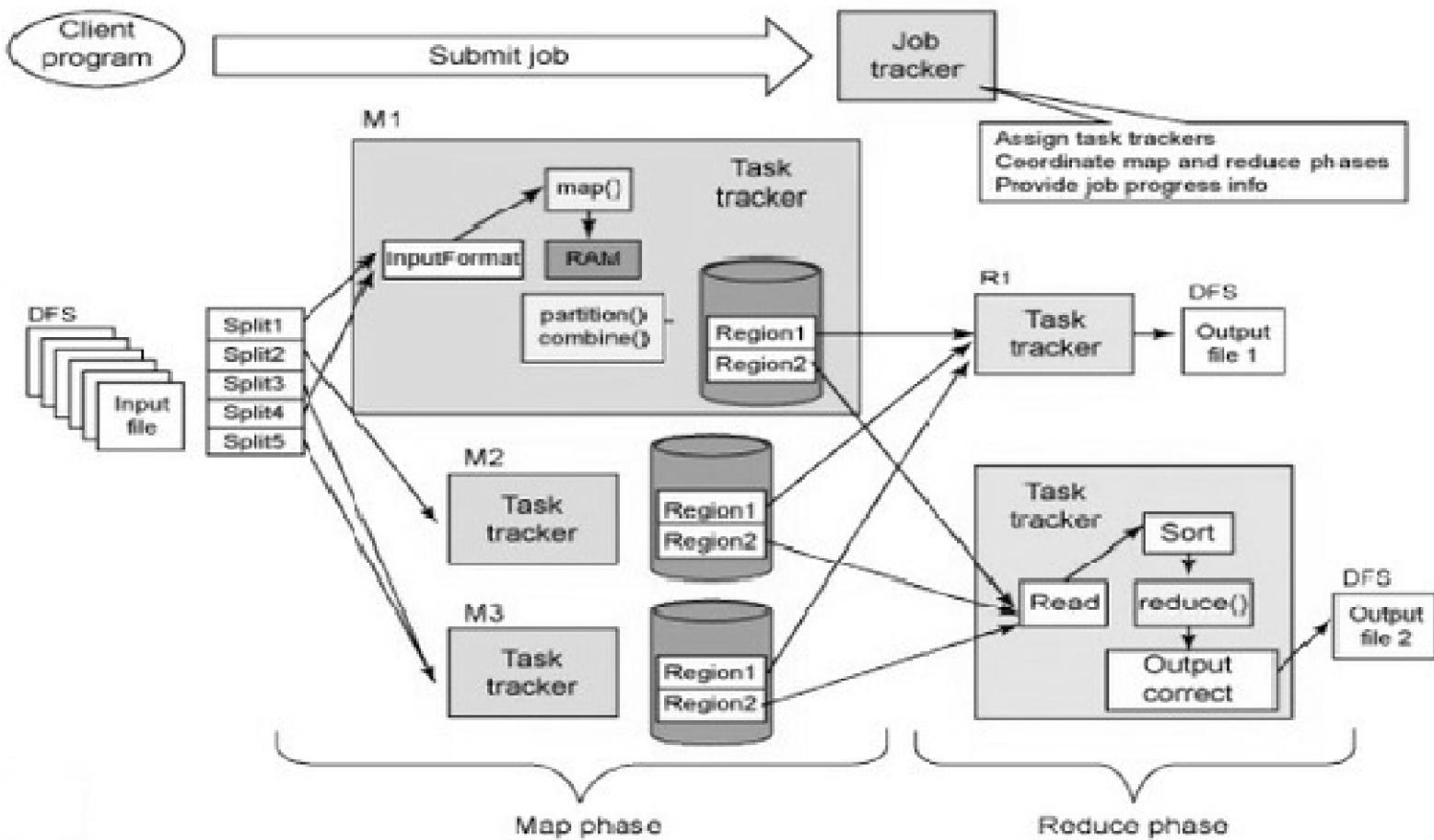
execuția unui job în map reduce

## Map Reduce - exemplu



$\langle k_1, v_1 \rangle \rightarrow$  transformare  $\rightarrow \langle k_2, v_2 \rangle \rightarrow$  reducere  $\rightarrow \langle k_3, v_3 \rangle$

# Map Reduce



Paşii

# **Securitatea Hadoop**

subțirică rău

# **Măsuri minime securizare hadoop**

ca de obicei!

# **Hadoop sau ElasticSearch?**

**Hadoop**

# **Exemplu aplicare Map Reduce**

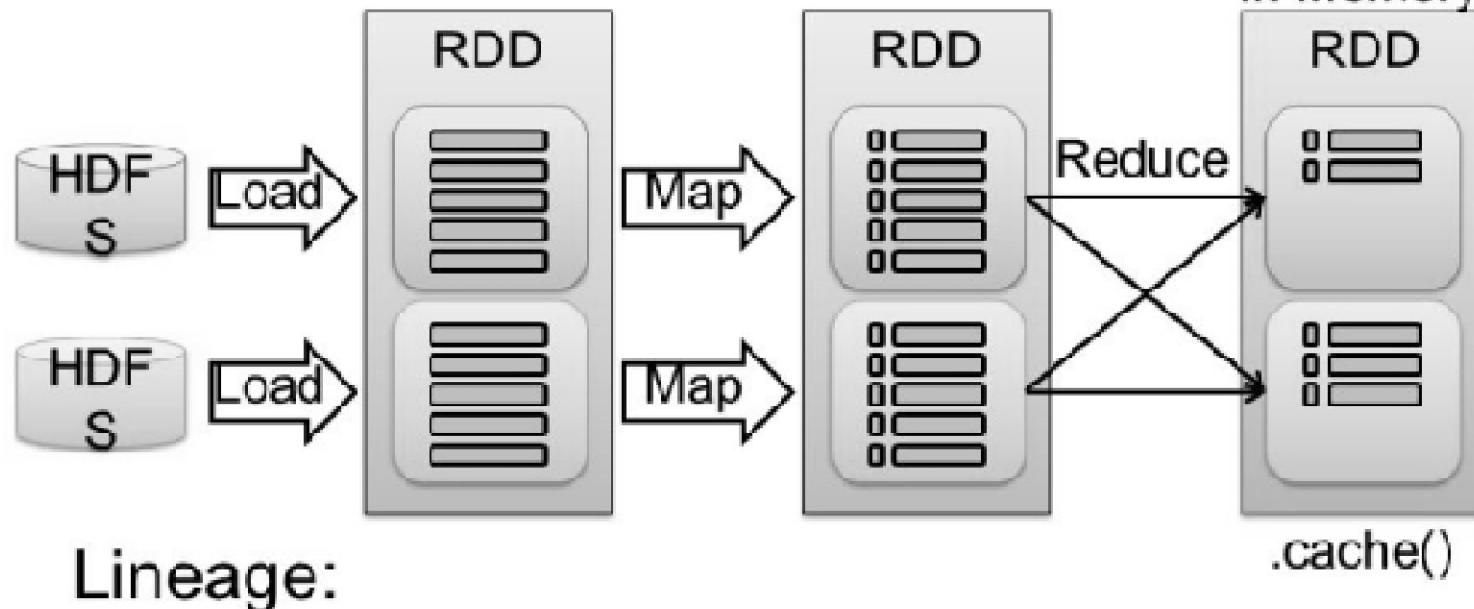
cel mai cunoscut!

# Spark

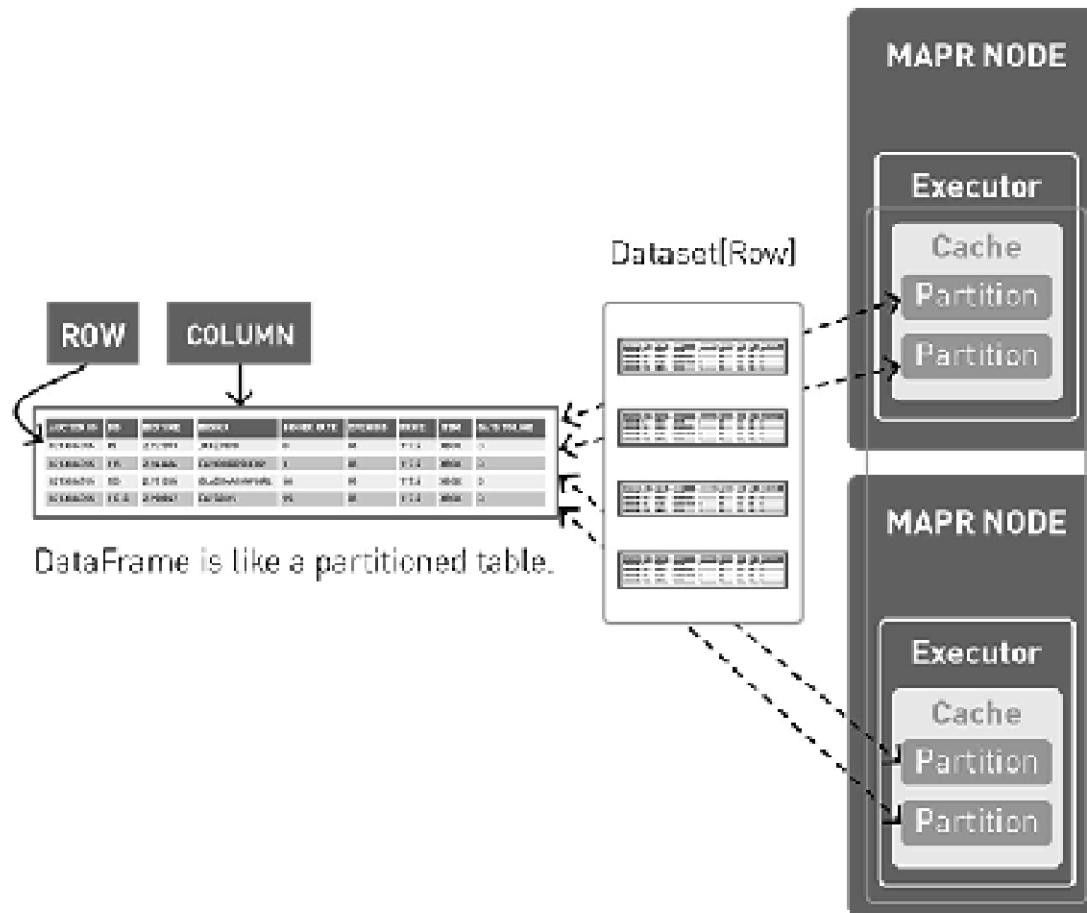
Zanaria et al., NSDI'12

## Resilient Distributed Datasets (RDD):

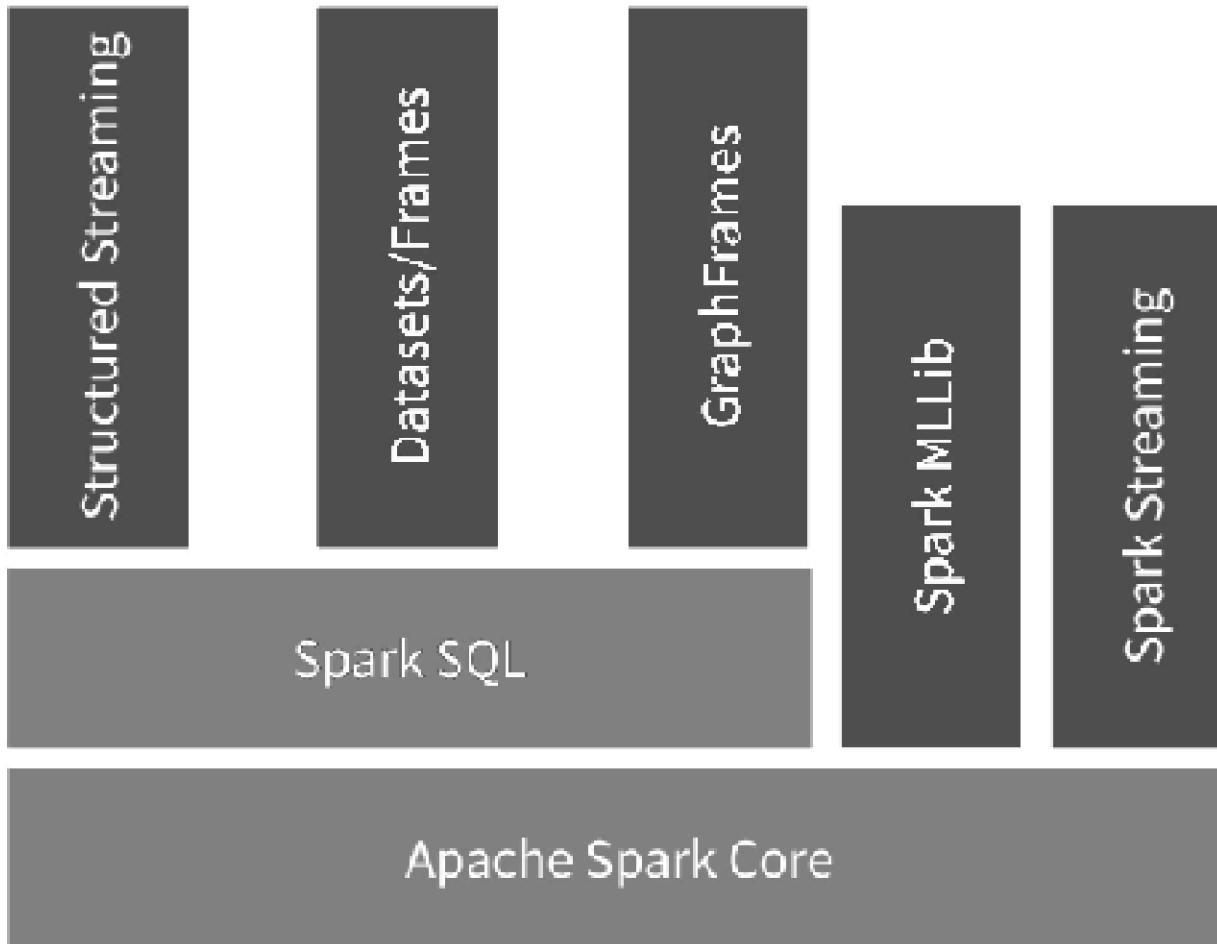
Persist  
in Memory



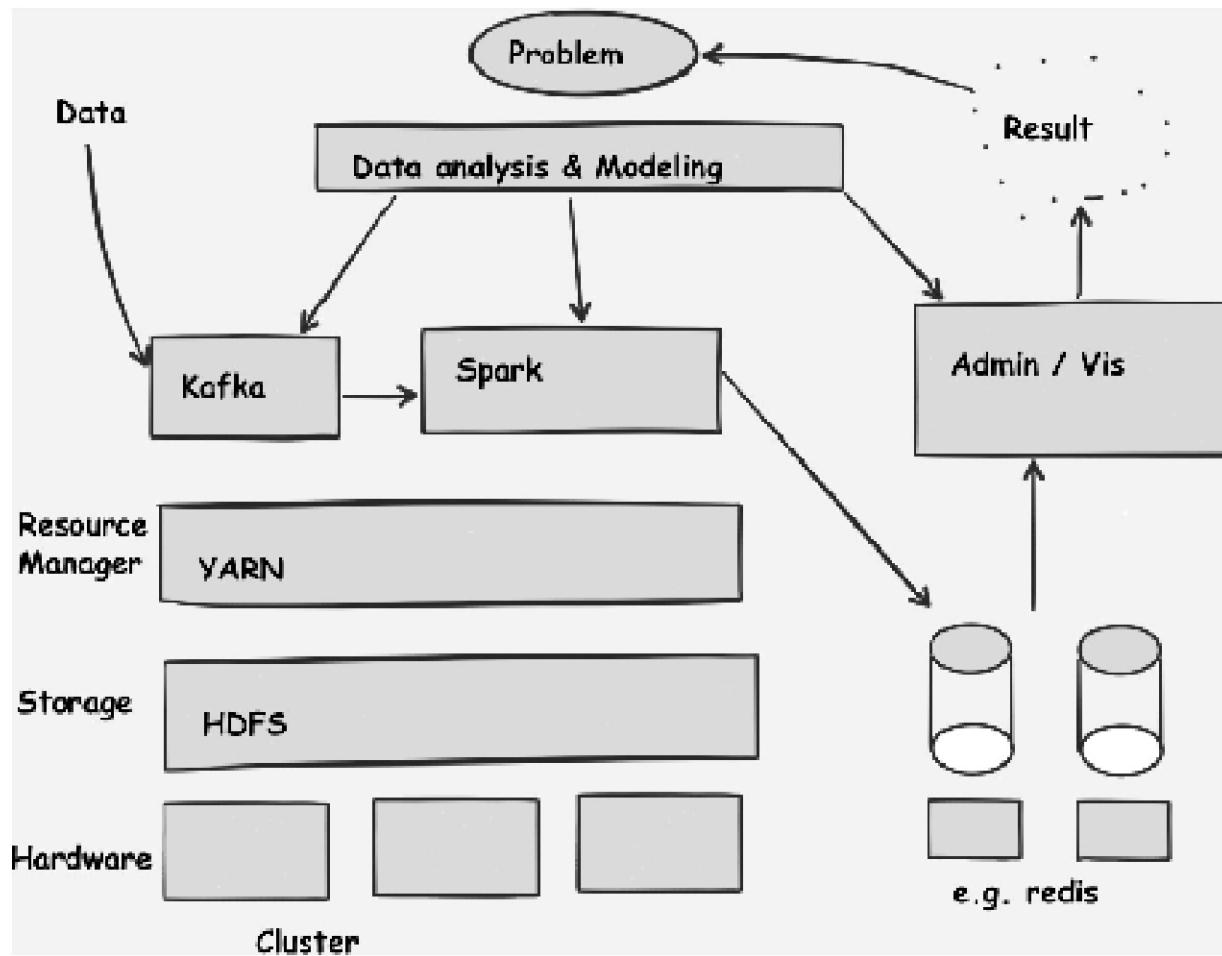
# Spark SQL



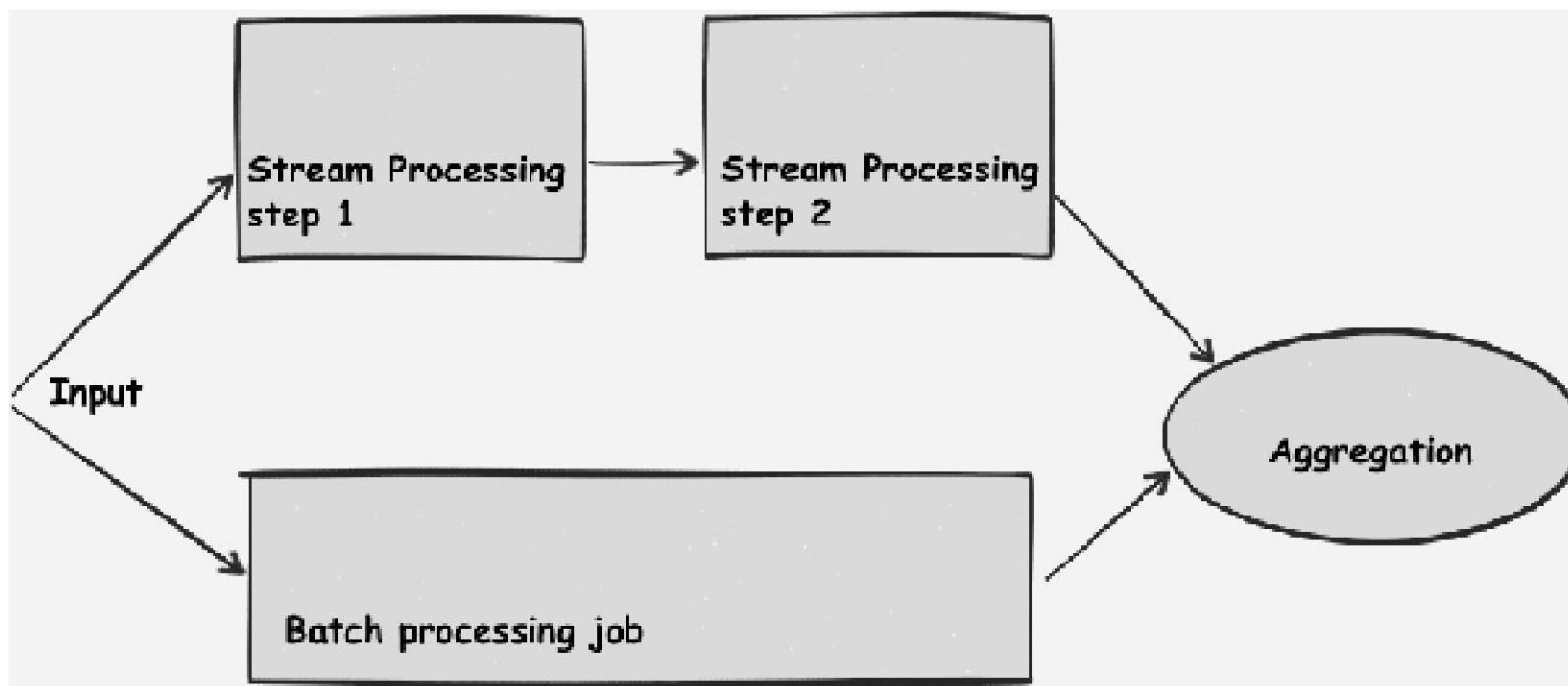
# Structura Spark



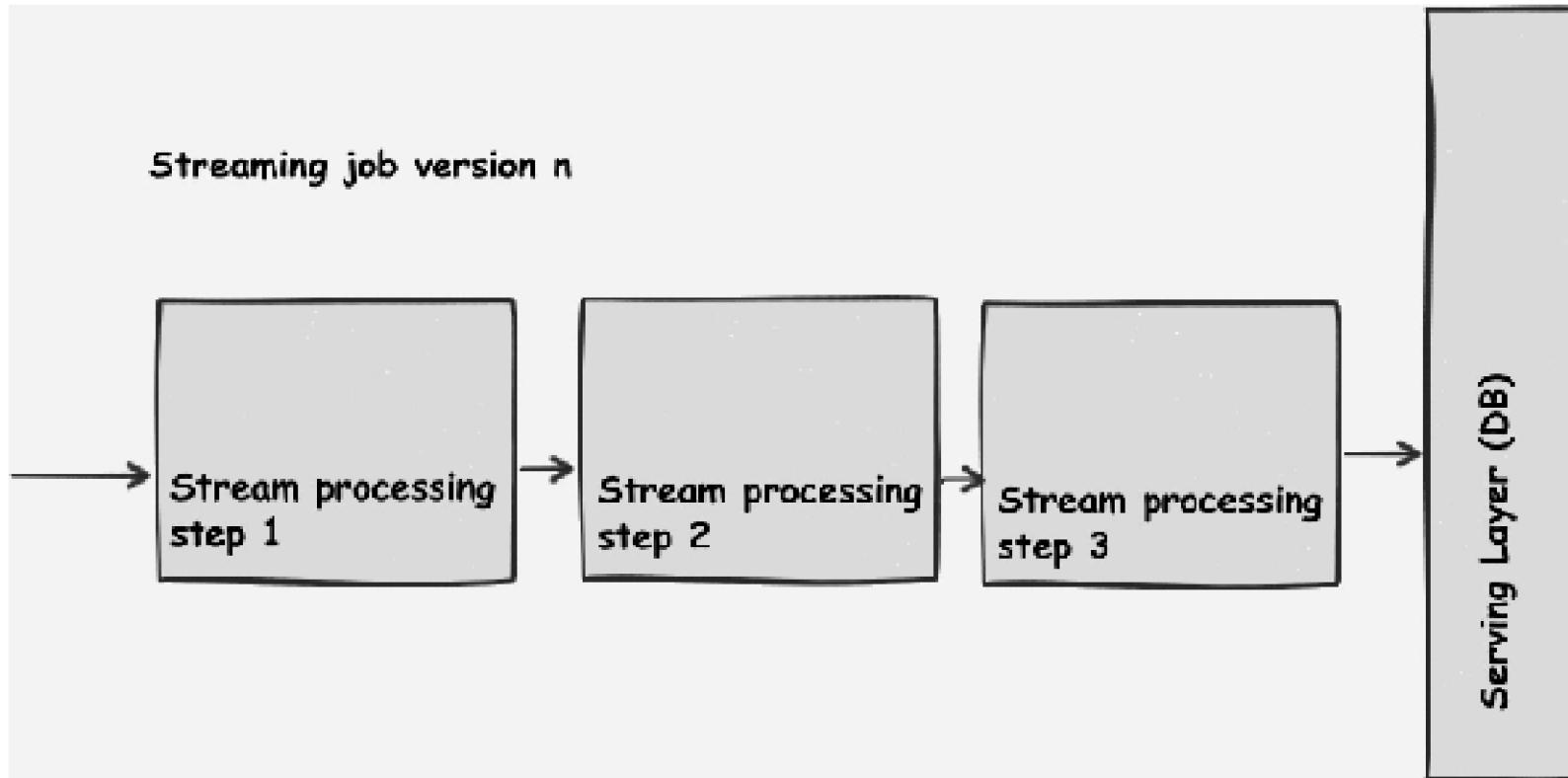
# O posibilă platformă de date bazată pe Spark



# Arhitectura Lambda

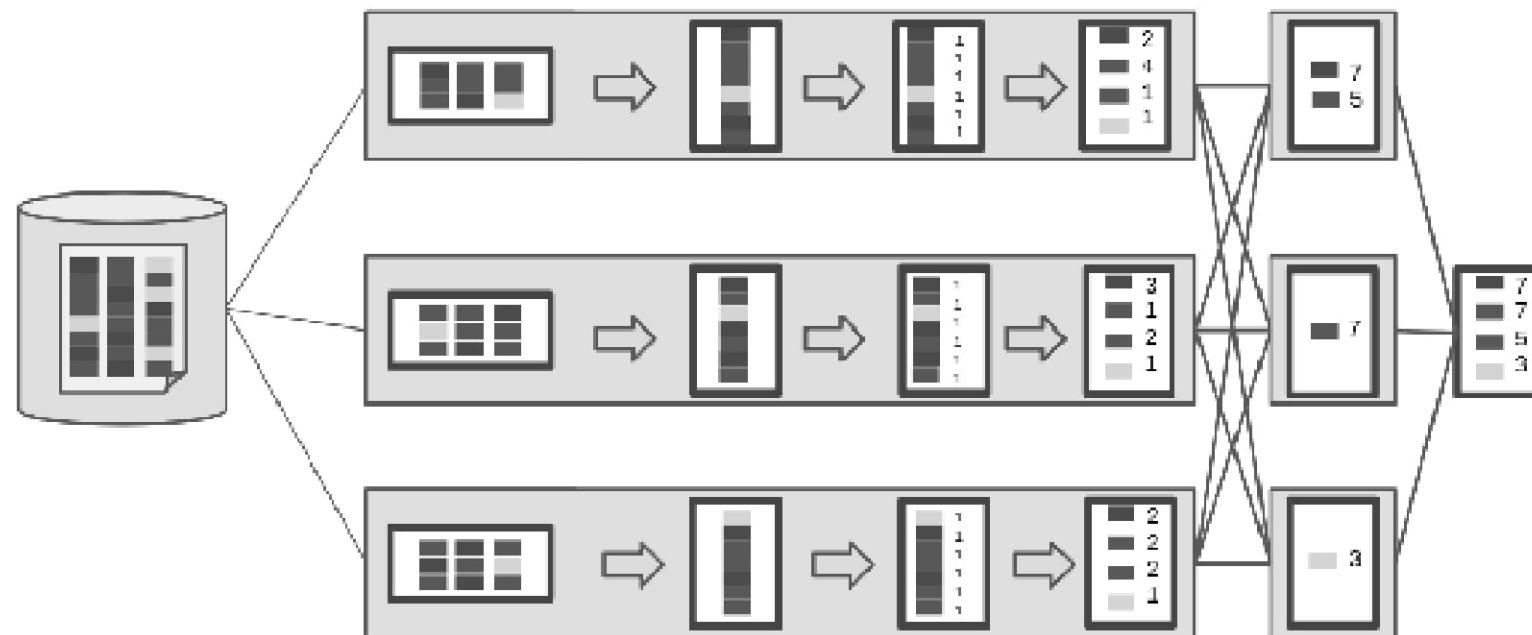


# Arhitectura Kappa



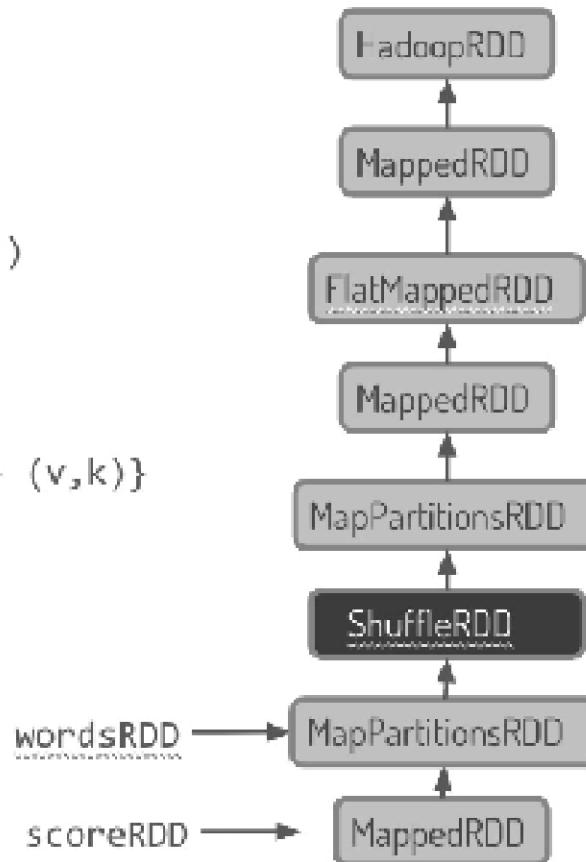
# RDD

```
.textFile("...") .flatMap(l => l.split(" ")) .map(w => (w,1)) .reduceByKey(_ + _)
```

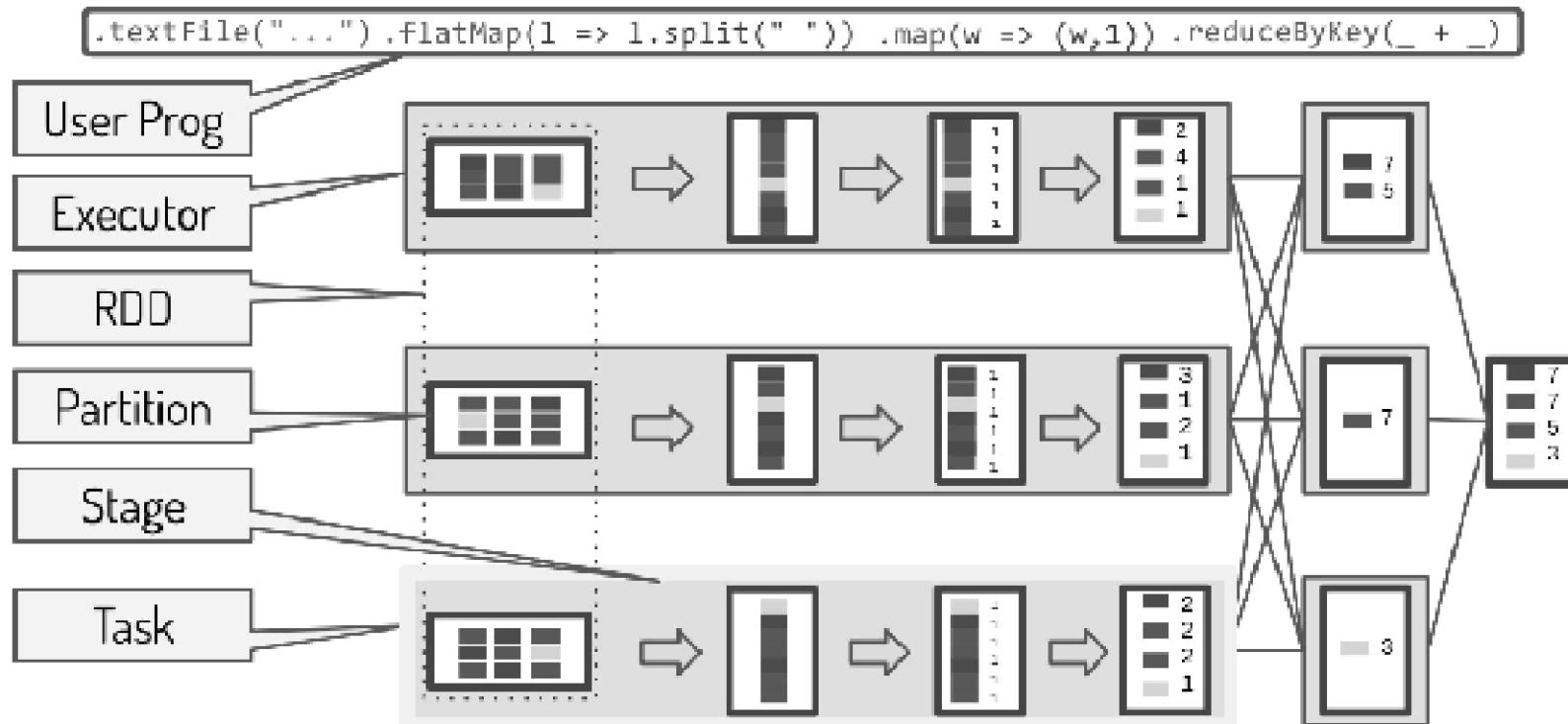


# RDD & DAG

```
val file = spark.textFile("hdfs://...")  
val wordsRDD = file.flatMap(line =>  
    line.split(" ")).  
    .map(word => (word, 1))  
.reduceByKey(_ + _)  
val scoreRDD = words.map{case (k,v) => (v,k)}
```



# Spark Components



- X