

Big Data Architektúrák és Elemző módszerek

Gombos Gergő, Laki Sándor

források: cognitiveclass.ai

Elérhetőségek

- Előadók:
 - Gombos Gergő
 - Laki Sándor
- honlap:
 - <http://ggombos.web.elte.hu>
 - <http://lakis.web.elte.hu>
 - (<http://bigdata.elte.hu>)
- email:
 - ggombos@inf.elte.hu
 - lakis@inf.elte.hu

A kurzusról

- Big Data rendszerek, architektúrák megismerése, alap szintű használata
 - Hadoop, HDFS, MapReduce
 - Spark, GraphX, Stream
- Data Science módszerek megismerése, algoritmusok használata
 - pandas, dataframe, numpy, jupyter
 - döntési fák, osztályozás, klaszterezés, dimenzió csökkentés
 - vizualizáció, matplotlib

A követelményekről

- Előadás:
 - Vizsga a félév végén
- Gyakorlat:
 - 1 beadandó: Hadoop MapReduce
 - (Házifeladatok)
 - 2 ZH a félév során:
 - Spark batch elemzés
 - Adatelemzés, vizualizáció

Big Data

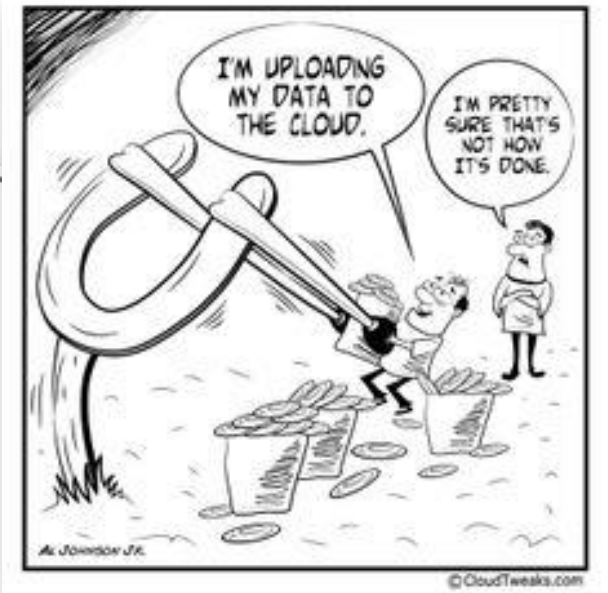


Big Data is like teenage sex:
everyone talks about it, nobody
really knows how to do it, everyone
thinks everyone else is doing it, so
everyone claims they are doing it.

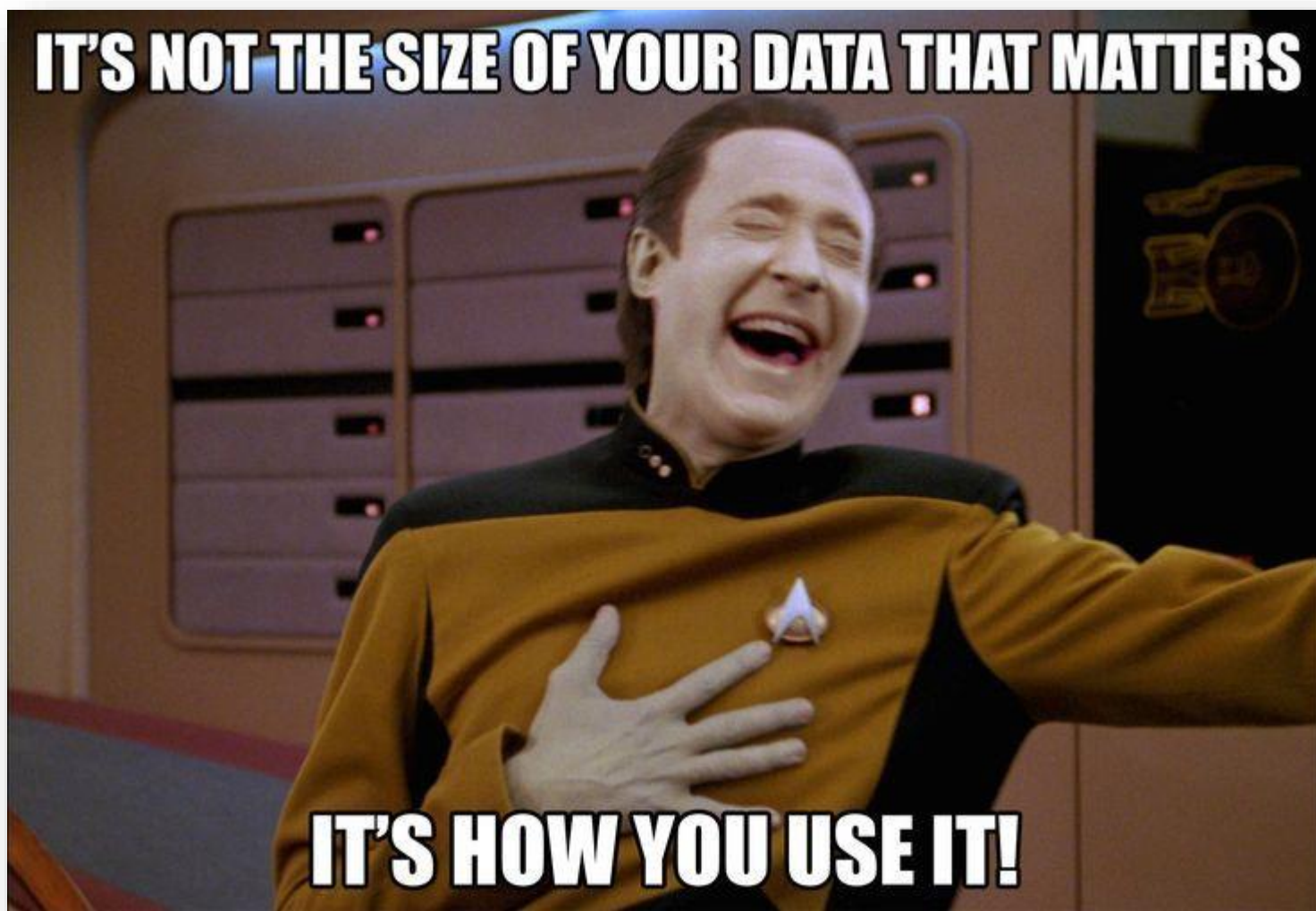
— *Dan Ariely* —

AZ QUOTES

A világ megváltozott!



Nem a méret a lényeg!



Mi az a Big Data?

BIG DATA DEFINITION



“Data that exceeds the processing capacity of conventional database systems. The data is too big, moves too fast, or doesn't fit the structures of your database architectures.”

Forbes

Source: Edd Dumbill, Forbes

40 ZETTABYTES

[43 TRILLION GIGABYTES]
of data will be created by 2020, an increase of 300 times from 2005

6 BILLION PEOPLE
have cell phones



WORLD POPULATION: 7 BILLION



Volume
SCALE OF DATA

It's estimated that **2.5 QUINTILLION BYTES** [2.3 TRILLION GIGABYTES] of data are created each day

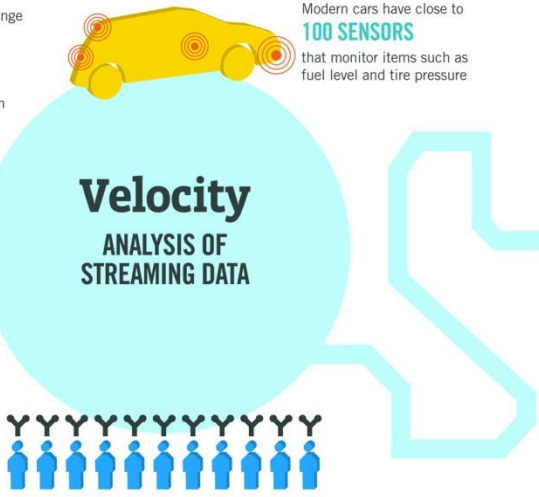


Most companies in the U.S. have at least **100 TERABYTES** [100,000 GIGABYTES] of data stored

The New York Stock Exchange captures **1 TB OF TRADE INFORMATION** during each trading session



By 2016, it is projected there will be **18.9 BILLION NETWORK CONNECTIONS** - almost 2.5 connections per person on earth



Velocity
ANALYSIS OF STREAMING DATA

Modern cars have close to **100 SENSORS** that monitor items such as fuel level and tire pressure



The FOUR V's of Big Data

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: **Volume, Velocity, Variety and Veracity**

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015 **4.4 MILLION IT JOBS** will be created globally to support big data, with 1.9 million in the United States



As of 2011, the global size of data in healthcare was estimated to be

150 EXABYTES [161 BILLION GIGABYTES]



30 BILLION PIECES OF CONTENT are shared on Facebook every month



Variety
DIFFERENT FORMS OF DATA

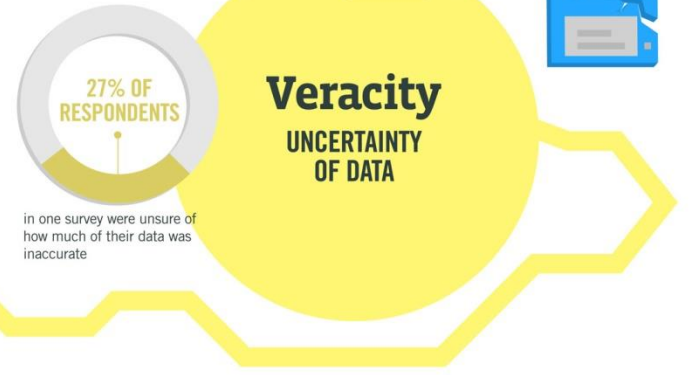
By 2014, it's anticipated there will be **420 MILLION WEARABLE, WIRELESS HEALTH MONITORS**

4 BILLION+ HOURS OF VIDEO are watched on YouTube each month



400 MILLION TWEETS are sent per day by about 200 million monthly active users

1 IN 3 BUSINESS LEADERS don't trust the information they use to make decisions



Veracity
UNCERTAINTY OF DATA

27% OF RESPONDENTS in one survey were unsure of how much of their data was inaccurate

Poor data quality costs the US economy around **\$3.1 TRILLION A YEAR**



Sources: McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTec, QAS



Volume - 640K Big Data?

1981



640K ought to be enough for anybody.

(Bill Gates)

Volume

Data Footprint of Humans



44

zettabytes

Projected volume of
global IT traffic by
2020



40

zettabytes

Volume of data
created by 2020, up
300% from 2015



2.3

zettabytes

Volume of data that
humans produce
every day

Source: IBM, Grazziti

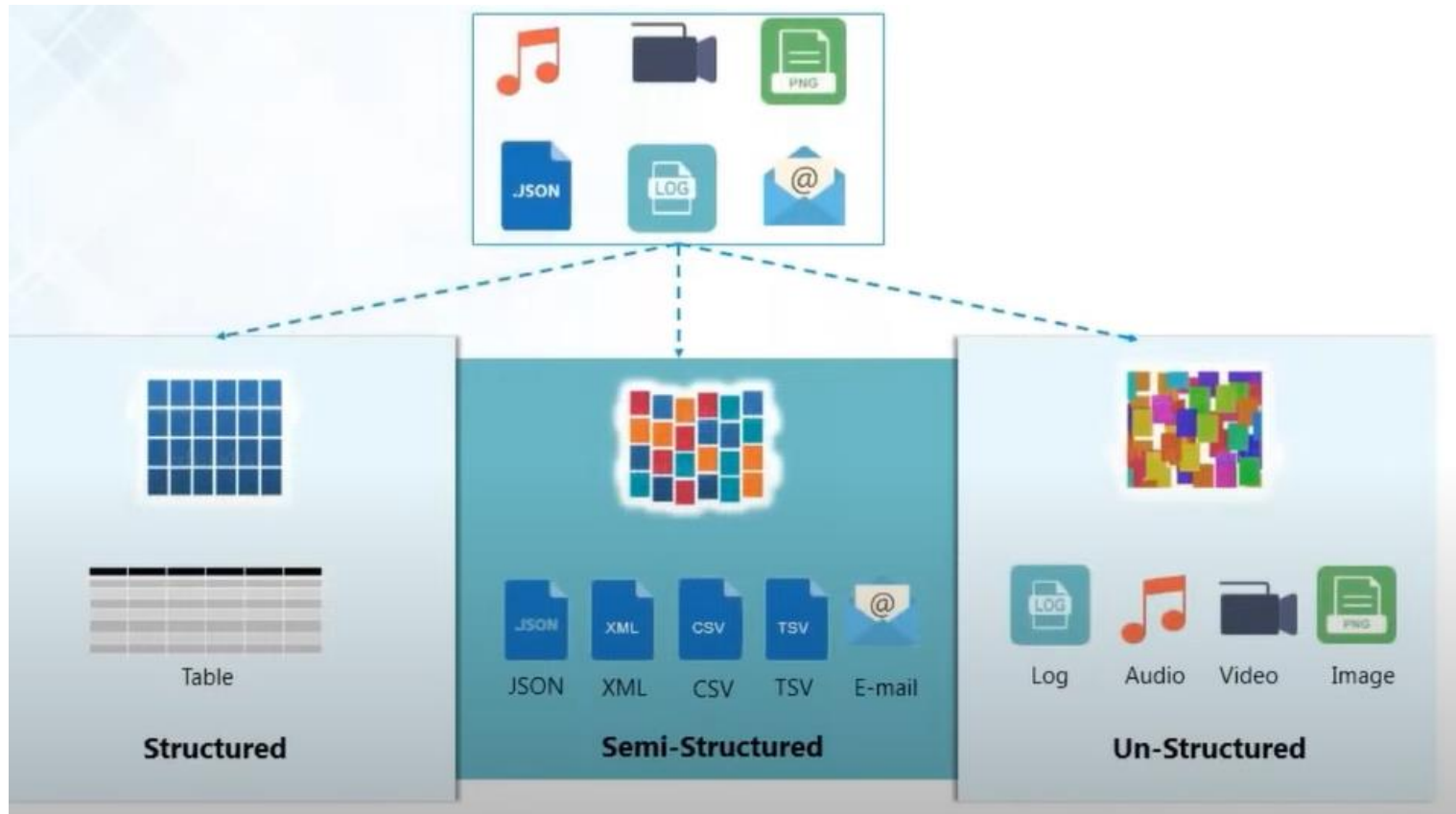
Velocity

THE INTERNET IN **2023** EVERY MINUTE



Created by: eDiscovery Today & LTMG

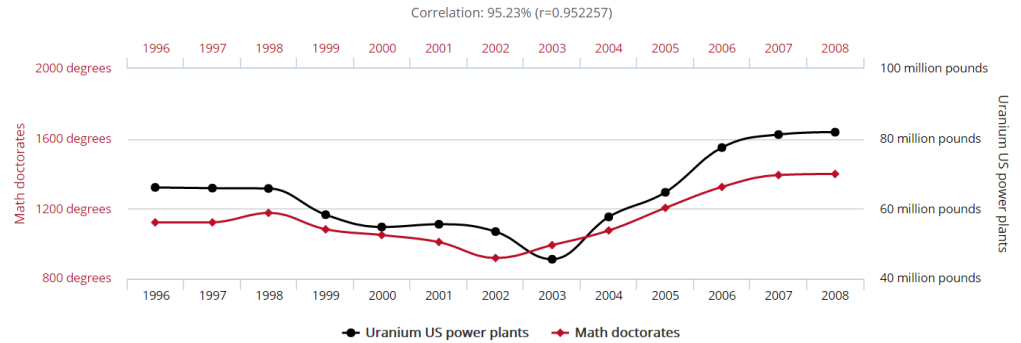
Variety (types of Big Data)



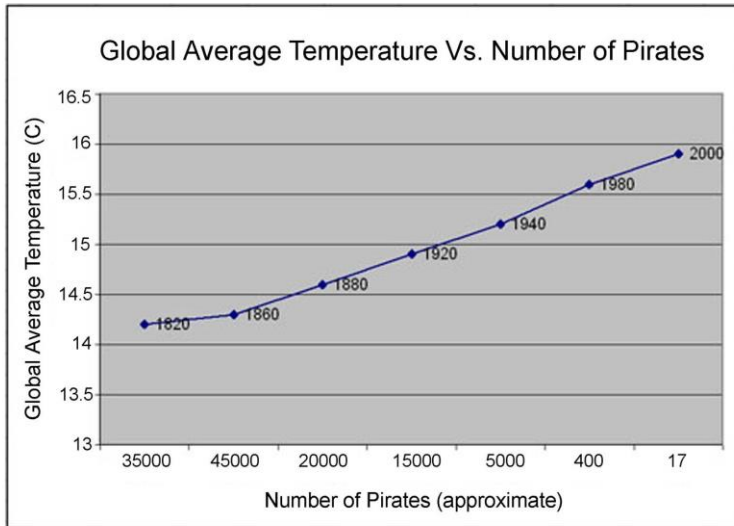
Veracity

Tényleg ez az összefüggés?

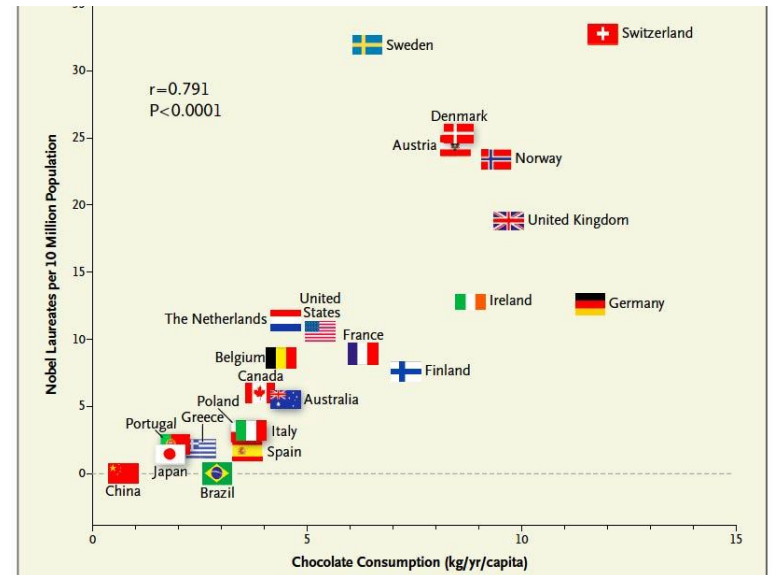
Math doctorates awarded
correlates with
Uranium stored at US nuclear power plants



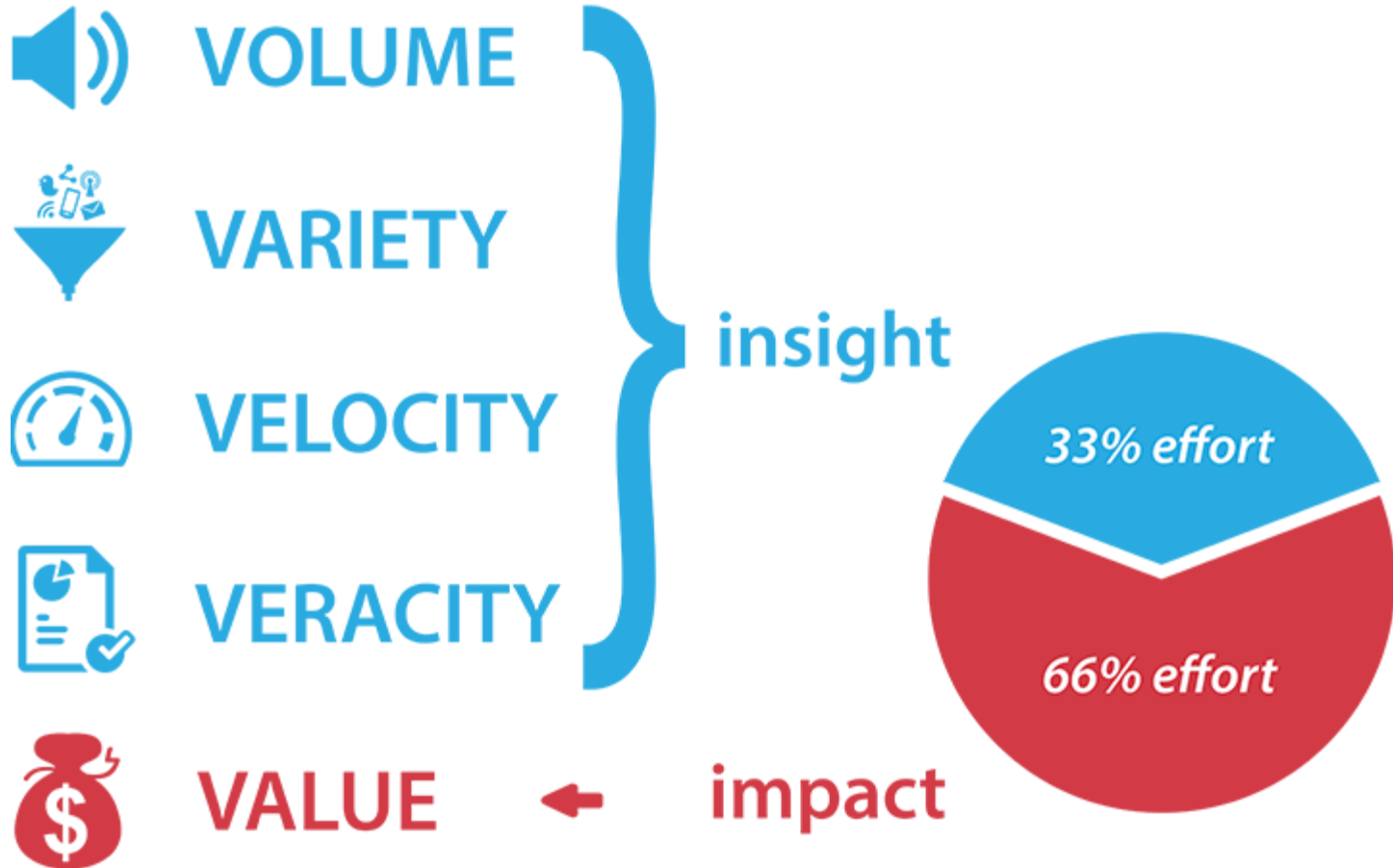
STOP GLOBAL WARMING: BECOME A PIRATE



WWW.VENGANZA.ORG



Big Data 5V



Netflix innovation relies on



Personalized
Video Ranking



Video-Video
Similarity Ranker



Project Cost
Predicting Algorithm



Artwork Visual
Analysis



Continue
Watching Ranker



Trending
Now Ranker

The impact of Big Data on businesses and people

Virtual personal assistants



Google Now

Siri knows what users mean when they ask her questions, she knows where they are, what time they are talking about and can use this information to look for restaurants of a particular type of food and check whether there reservations are available.

Google Now makes recommendations before users ask for them, especially when it is linked up to the user's calendar and location sensing is enabled on the user's phone. Google Now knows where the user is and where he/she needs to be, it can tell users about things like traffic or the weather before you even ask for it.



The V's Evolution of Big Data



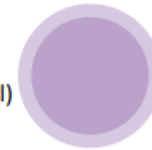
9V's (Owais,2016)

- Volume
- Variety
- Velocity
- Veracity
- Value
- Visualization
- Variability
- Validity
- Volatility



10V's (Data Science Central)

- Volume
- Variety
- Velocity
- Veracity
- Value
- Variability
- Validity
- Venue
- Vocabulary
- Vagueness



17V's (Panimalar,2017)

- Volume
- Variety
- Velocity
- Veracity
- Value
- Variability
- Validity
- Venue
- Vocabulary
- Vagueness
- Volatility
- Visualization
- Viscosity
- Virality
- Verbosity
- Voluntariness
- Versality

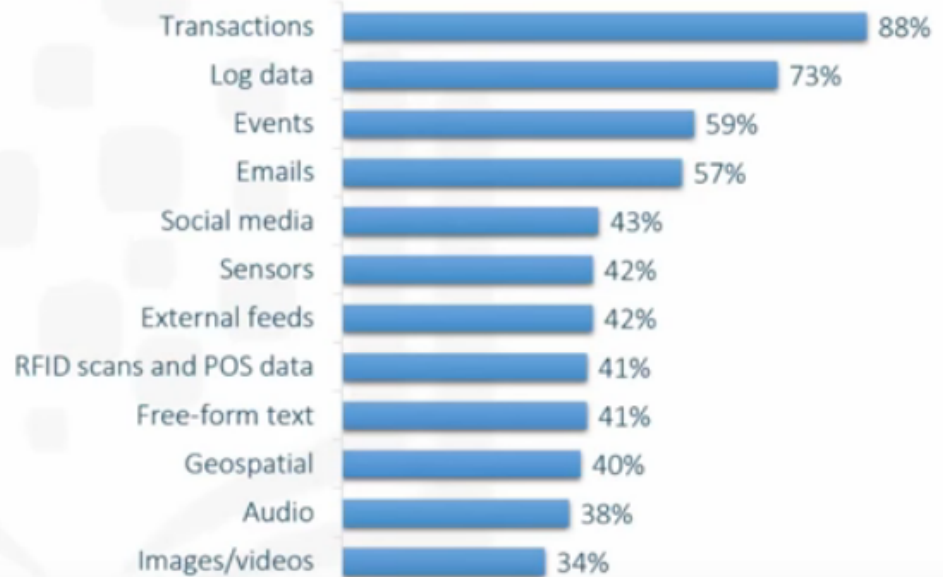
Where is all the data coming from?

There are three major sources of Big Data:

- People-generated data
- Machine-generated data
- Business-generated data

Big data sources

Multiple answers allowed

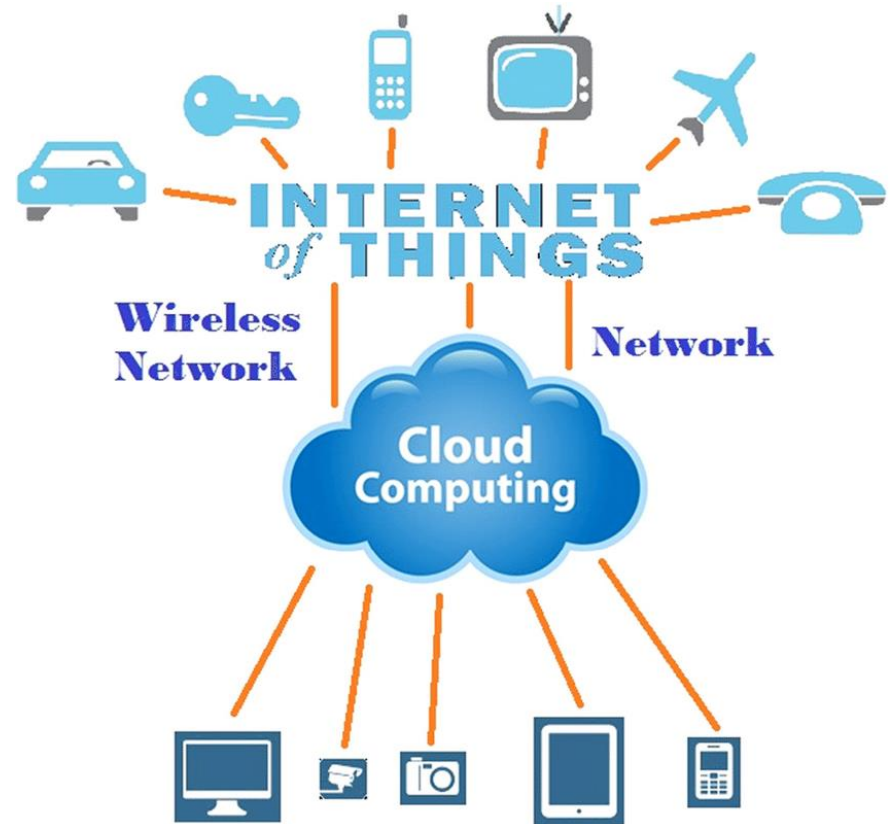
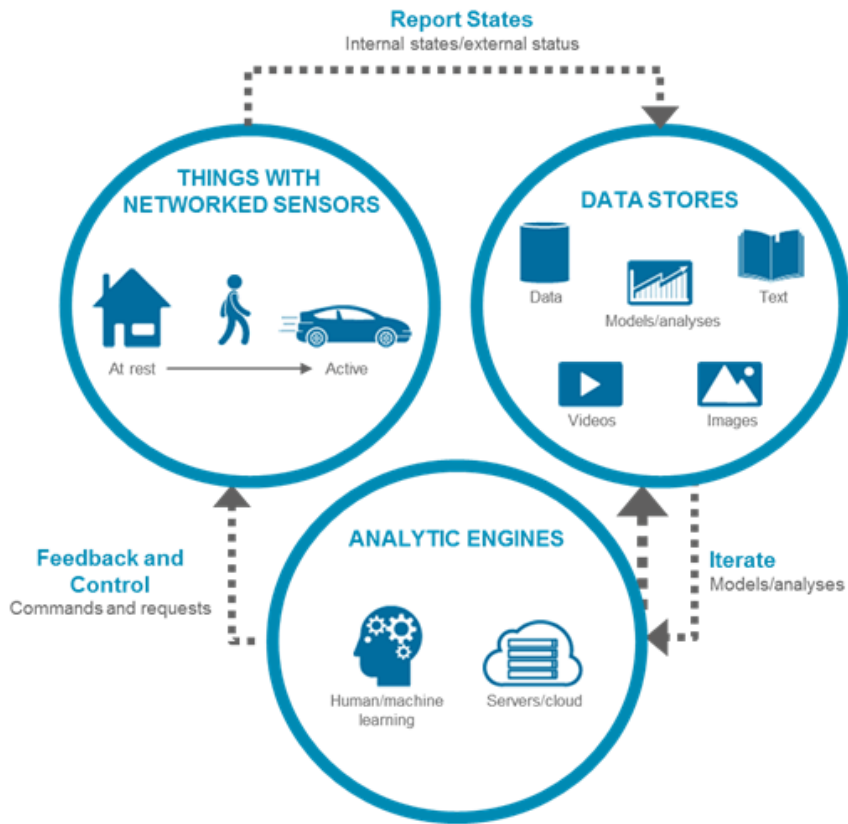


Source: IBM, Saïd School of Business, 2012

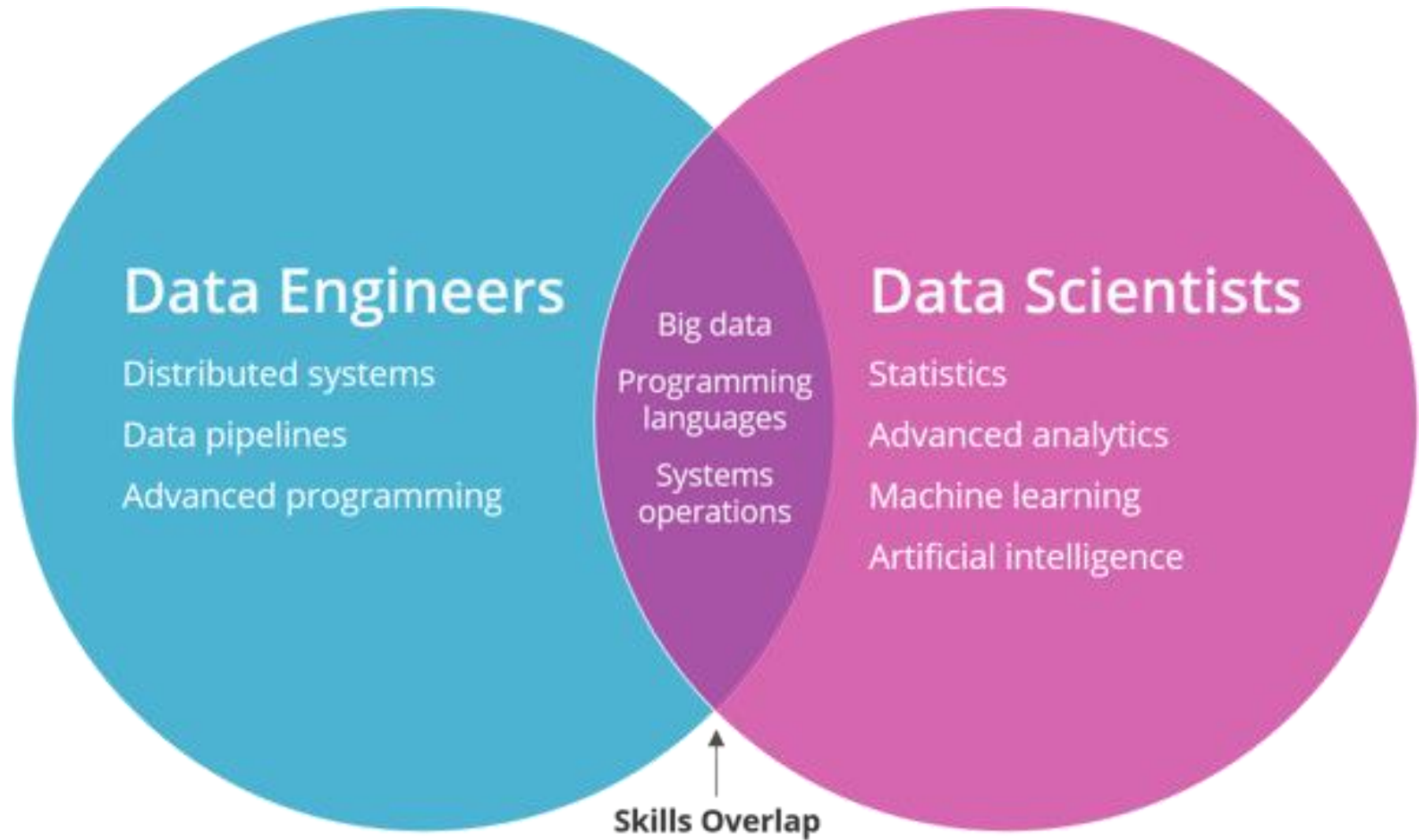


IoT – Internet of Things

Interaction Between the Three Components of the Internet of Things



Big Data Skillek



Data Engineers

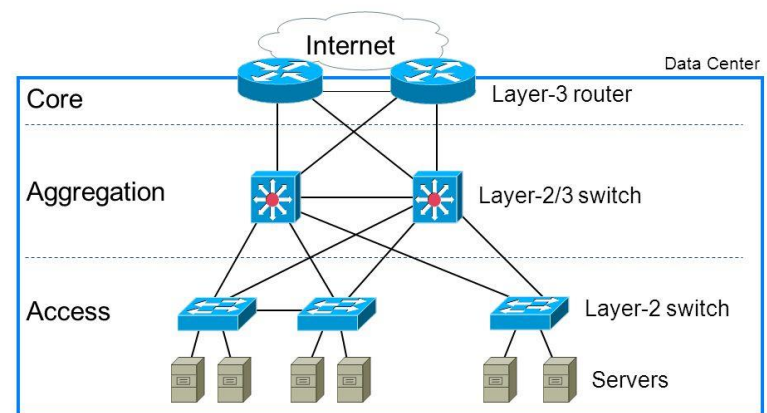
Motiváció – Hagyományos rendszerek

- Hagyományos rendszerek
 - Tipikusan egygépes rendszerek
 - Relatív „kis” méretű statikus adatok
 - Komplex feladatok elvégzése ezeken az adatokon működik
 - Mi történik ha nem fér el az adat? (Google példa)
 - 10 milliárd weblapot dolgoz fel naponta (átlagos méret 20 KB)
10 milliárd*20 kB = 200TB
 - Lemez olvasási sebesség: 50MB/s => 4 milliárd mp ≈ 46+ nap
- Számítási kapacitás növelése
 - Gyorsabb processzor
 - Több memória

Big Data Architektúrák

- Big Data Architektúra
 - Folyamatos adatfolyamok (napi > 1TB)
 - Több gép klaszterben
 - Az adatok elosztva tárolódnak (pl.: HDFS)
 - Számítások párhuzamosítva (MapReduce)
- Kapacitás növelés
 - Új gépek kapcsolása a klaszterhez

Common Data Center Topology



Technológiák (nem teljes)



Impala



Zookeeper



Pig

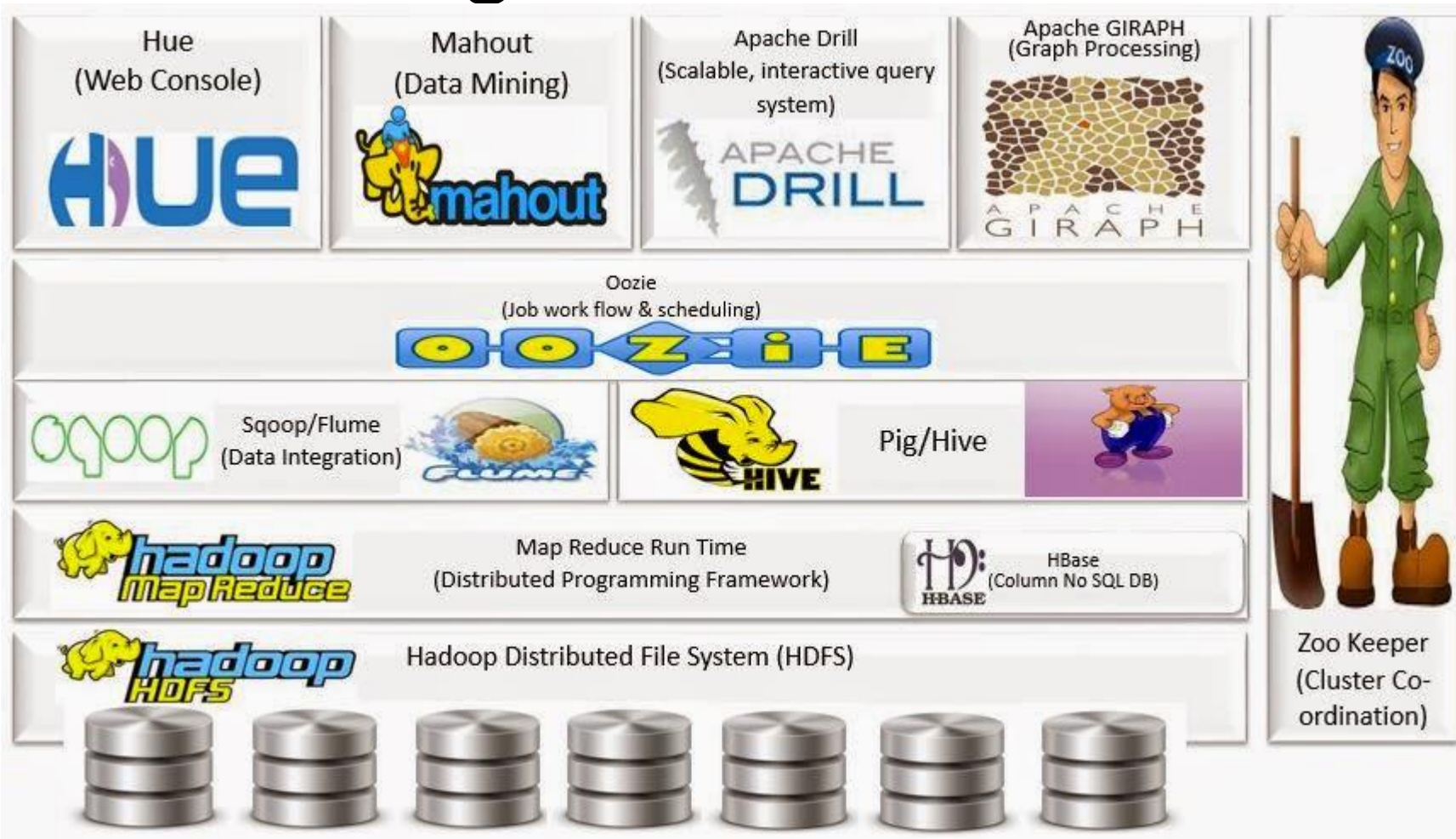


<https://hadooecosystemtable.github.io/>

Kik használják?



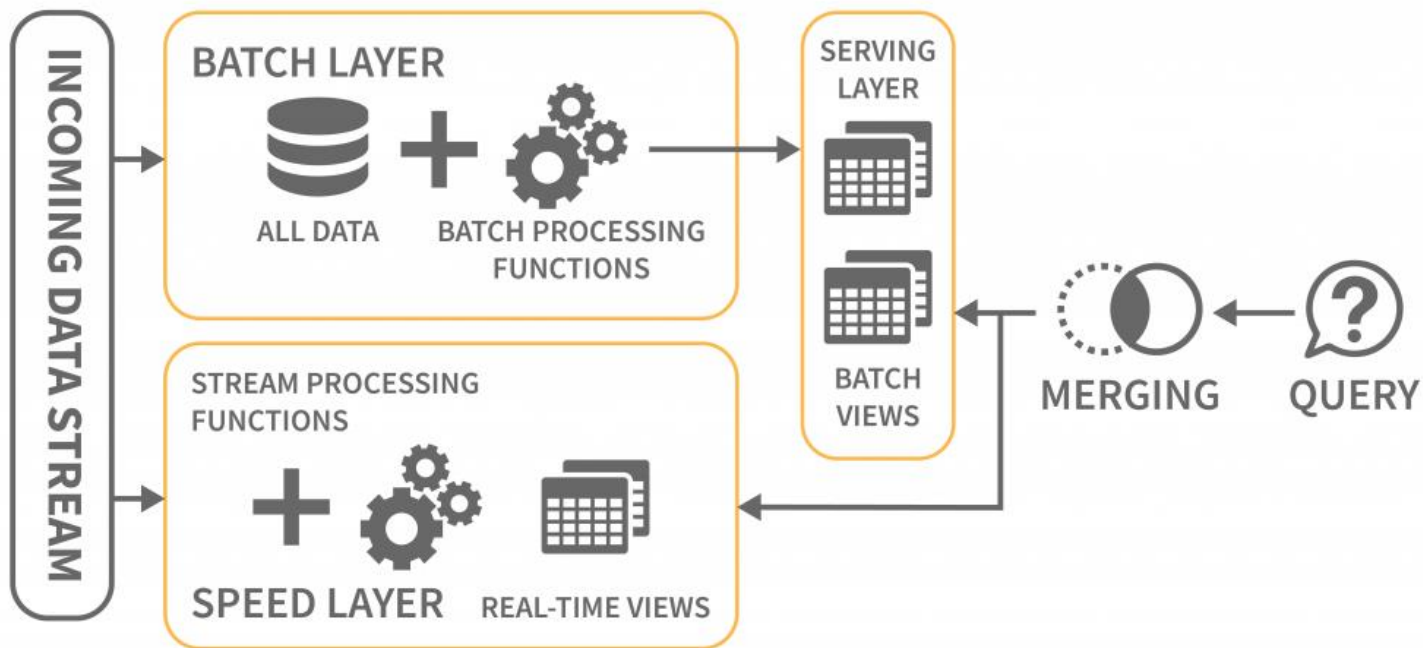
Példa Big Data Architektúrára



<https://hadooecosystemtable.github.io/>

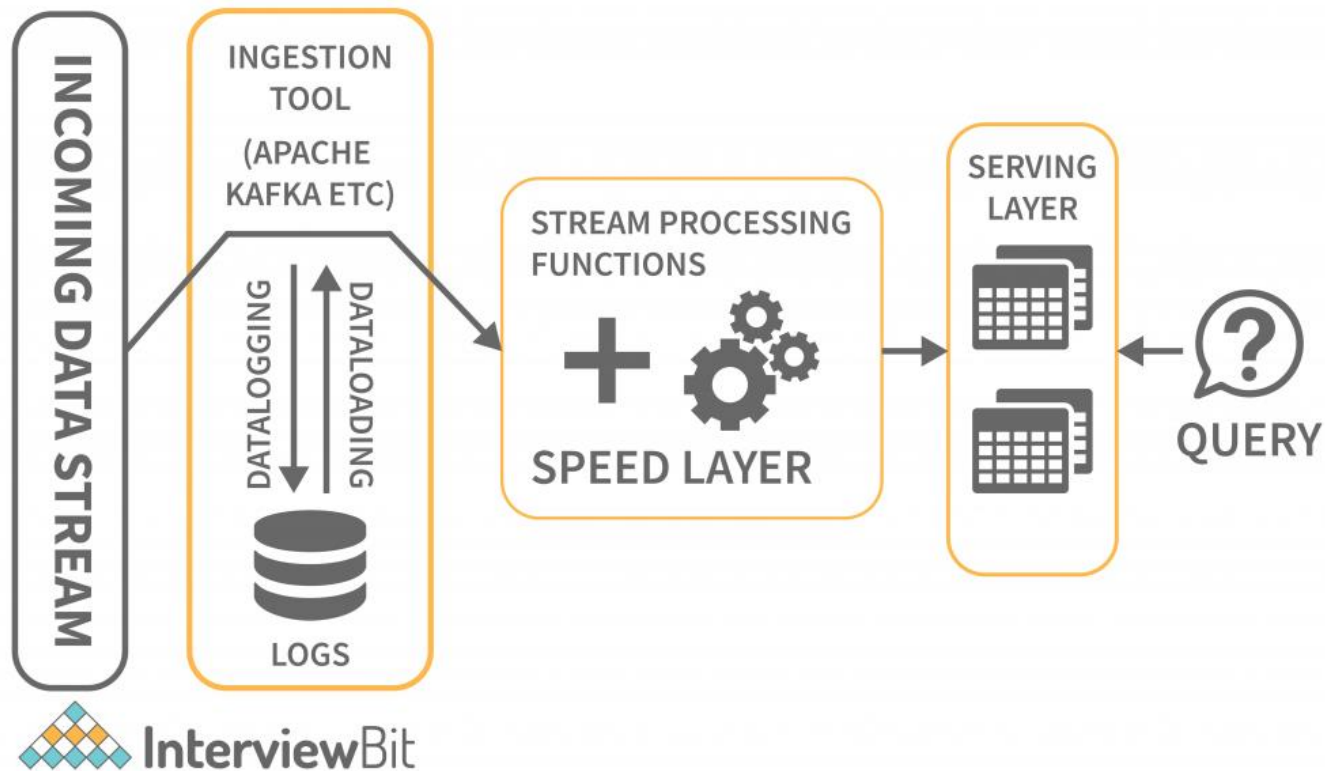
Big Data Architektúra

LAMBDA ARCHITECTURE



Big Data Architektúra

KAPPA ARCHITECTURE



Disztribúciók

- Apache
 - Az elemek egyesével érhetőek el.
 - Opensource
 - Telepítő: Ambari
- Cloudera, Hortonworks
 - Komplet telepítő rendszer
 - Support (fizetős)
 - Trainingek (fizetős)

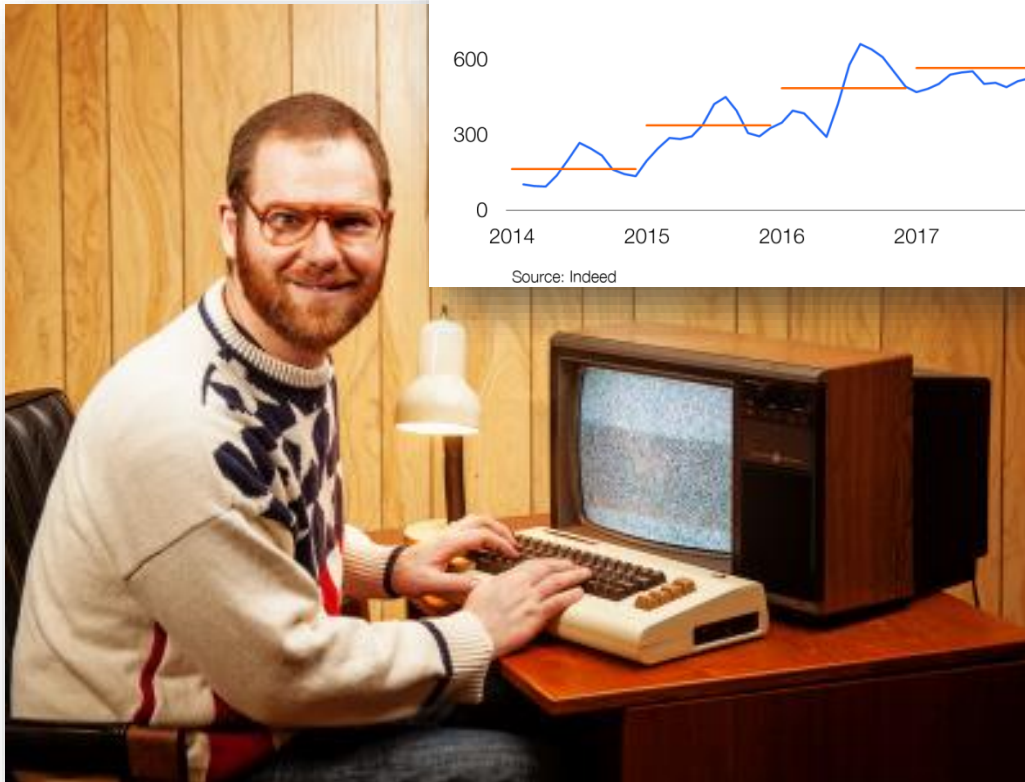
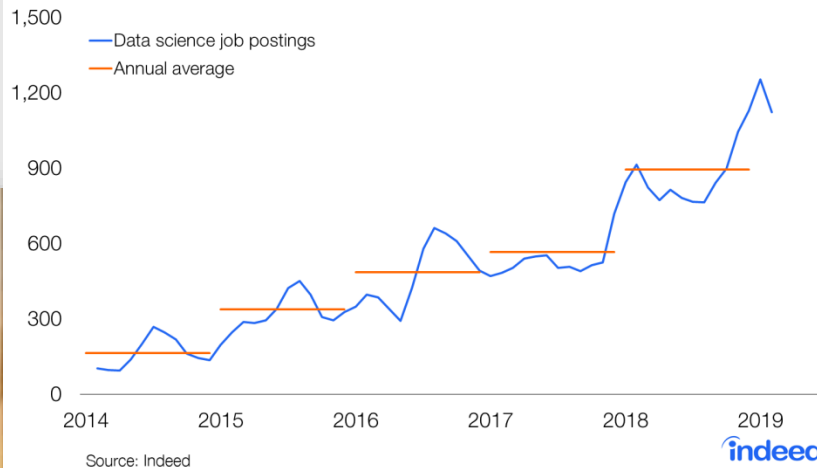
Data Scientist

Ki az a Data Scientist?

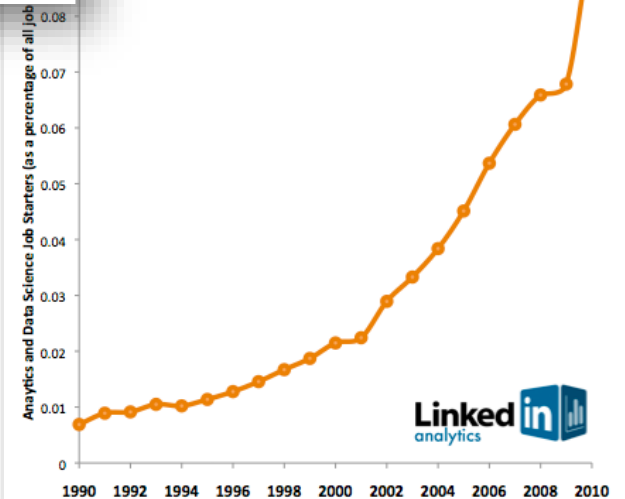


Data Scientist a „legszexibb” foglalkozás a 21. században!

In Australia, demand for data scientists is booming
Australian data science job postings, per million job postings, 3-month moving average



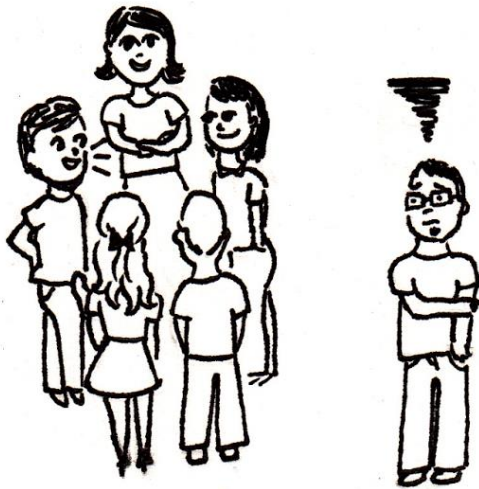
Analytics and Data Science Job Growth



Data Scientist

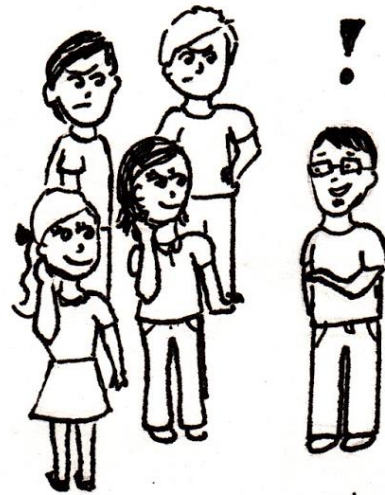
The Rise of Data Scientists

BEFORE



nobody cared for a
"math geek" in parties.

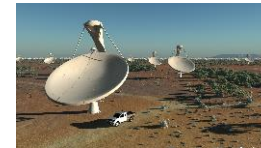
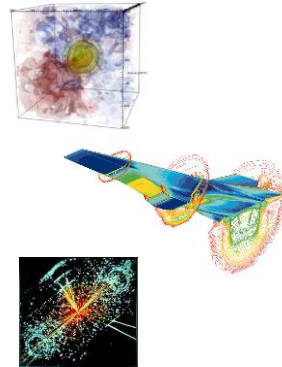
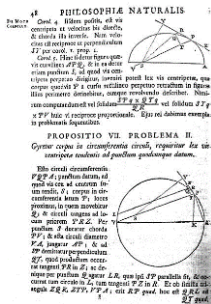
NOW



People love ~~math~~geeks
data scientists!

RK

Adat vezérelt világ!



4000 years

500 years

~50 years

Today

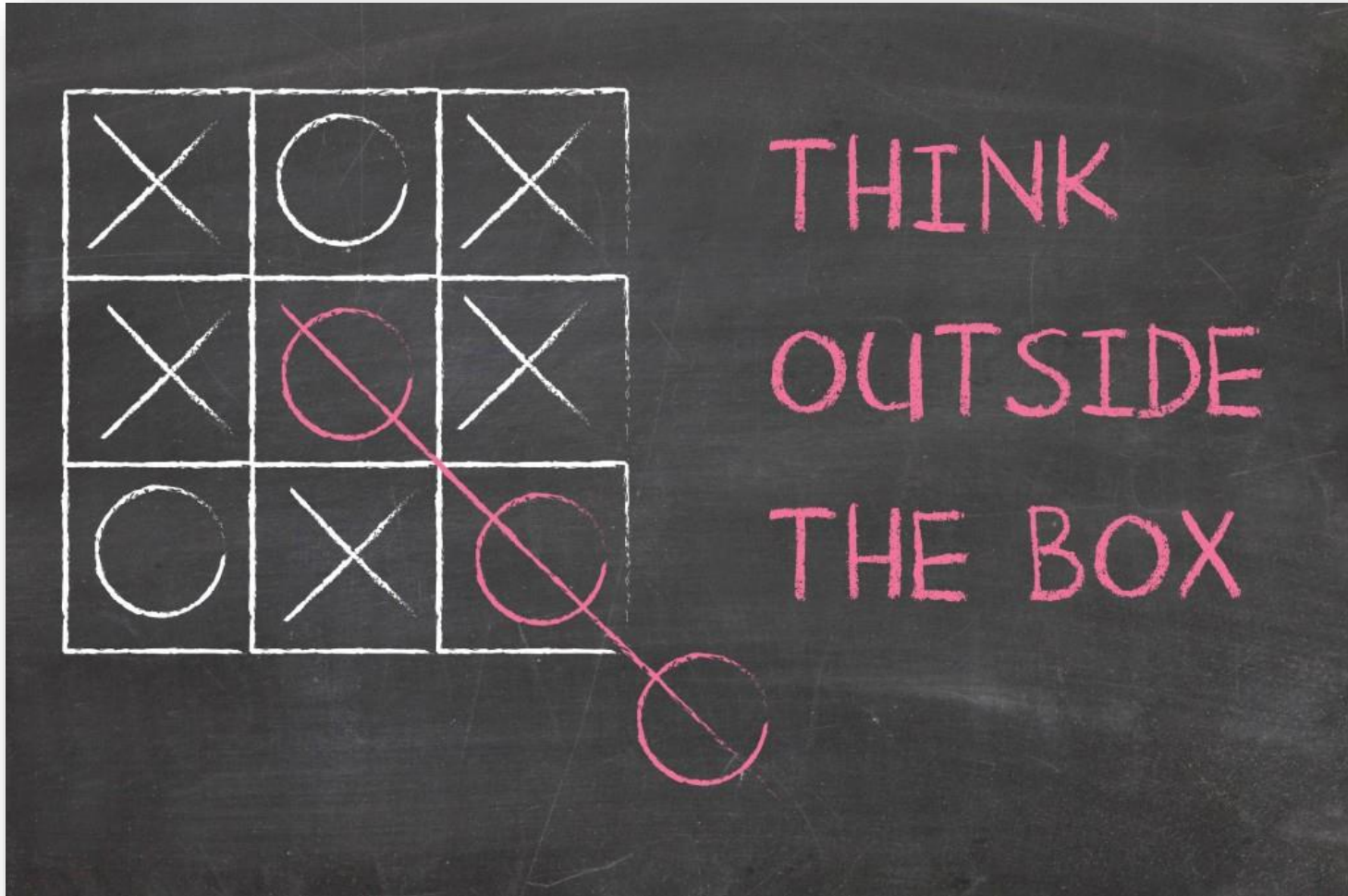
1 – Empirical observations

2 - Generalization Theoretical models

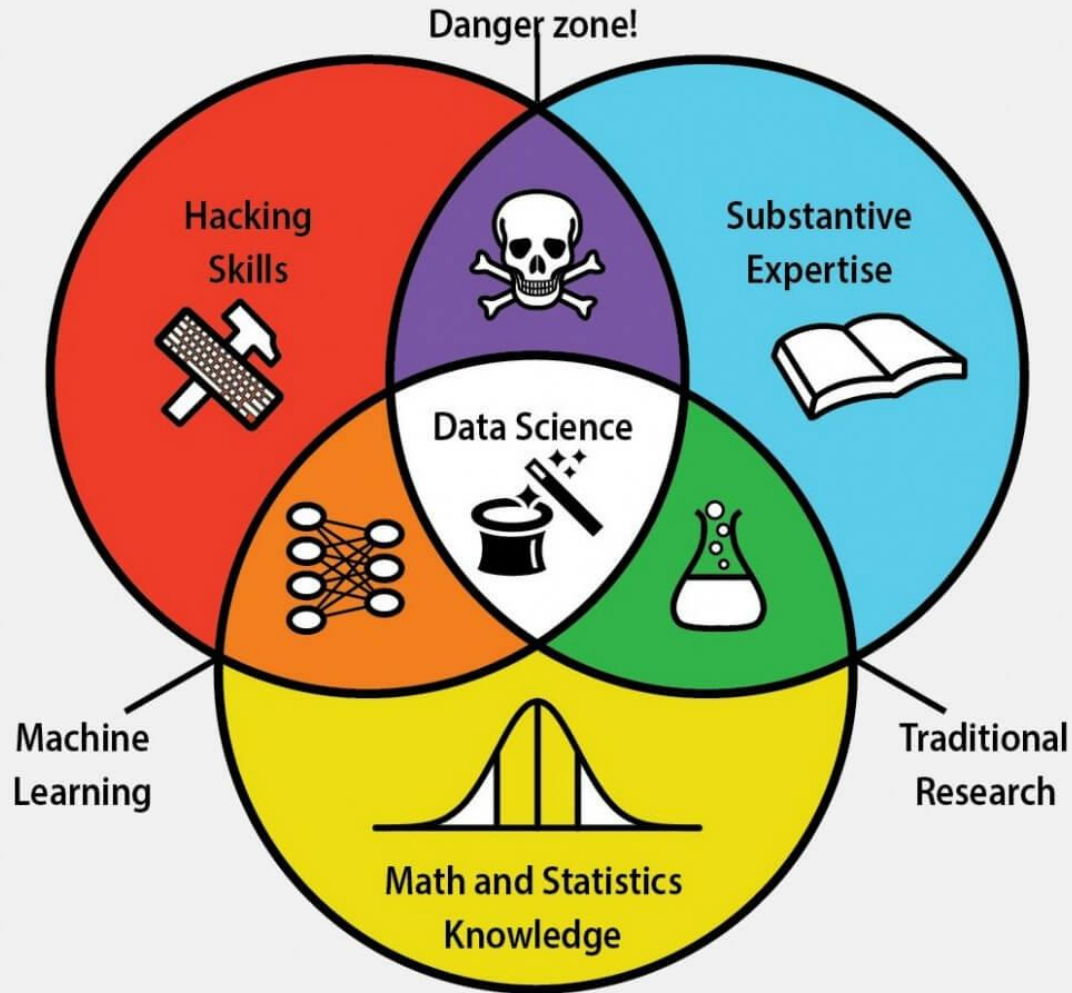
3 - Simulations Computational sciences

4 - Data-driven science eScience

Gondolkozz Másképp!



DATA SCIENCE SKILLSET



A Data Scientist munkája – a munkatársak szerint

$$\hat{\mathbf{r}} = \frac{x\hat{\mathbf{X}} + y\hat{\mathbf{Y}} + z\hat{\mathbf{Z}}}{\sqrt{x^2 + y^2 + z^2}}$$
$$\hat{\boldsymbol{\theta}} = \frac{(x\hat{\mathbf{X}} + y\hat{\mathbf{Y}})z - (x^2 + y^2)\hat{\mathbf{Z}}}{\sqrt{x^2 + y^2 + z^2}\sqrt{x^2 + y^2}}$$
$$\hat{\boldsymbol{\varphi}} = \frac{-y\hat{\mathbf{X}} + x\hat{\mathbf{Y}}}{\sqrt{x^2 + y^2}}$$

A Data Scientist munkája – az anyukája szerint



A Data Scientist munkája – valójában

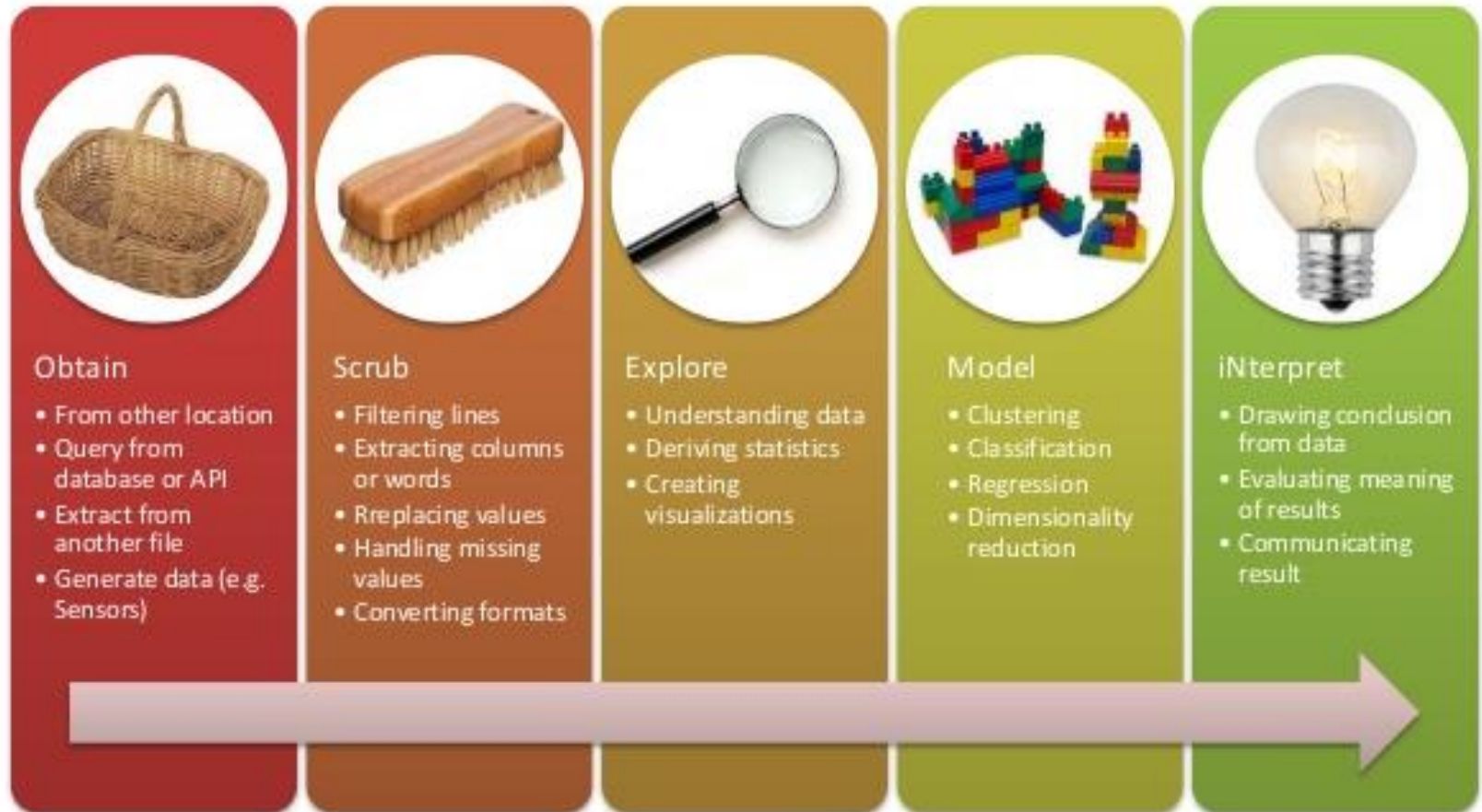
```
import pandas as pd
```

```
import numpy as np
```

```
SELECT * FROM Data WHERE...
```

@beckerfuffle #pyDataNYC

The Data Science is OSEMN



Source: [A Taxonomy of Data Science](#)

Big Data kihívások



hadoop történelem

- 2004 – Google publikálja a MapReduce technikát
- 2006 – Apache projekt lett Yahoo! támogatással
- További támogatók:

amazon.com[®]

IBM

last.fm

VISA

facebook

LinkedIn[®]

twitter

Google

JPMORGAN CHASE & CO.

The New York Times

YAHOO![®]

Kiknek jó a hadoop?

“... to create building blocks for programmers who just happen to have ***lots of data to store, lots of data to analyze, or lots of machines to coordinate***, and who don't have the time, the skill, or the inclination to become distributed systems experts to build the infrastructure to handle it.”

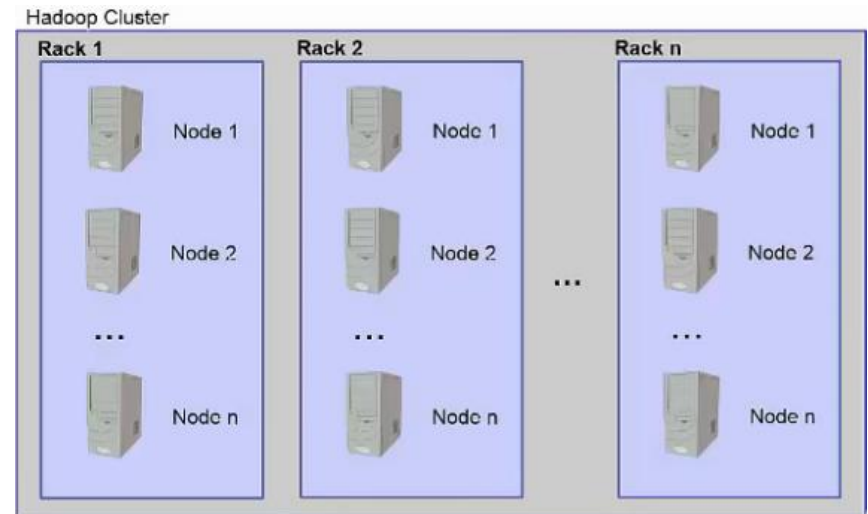
Tom White
Hadoop: The Definitive Guide

Kik használják?



Terminológia

- Klaszter:
 - Összetartozó számítógépek halmaza
- Rack:
 - Egy fizikai vagy virtuális csoportba tartozó gépek, amelyek jellemzően azonos áram és hálózati elérésen keresztül érhetőek el.
- Node:
 - Számító egységek, tipikusan a számítógépek



Hadoop Sorts a Petabyte in 16.25 Hours and a Terabyte in 62 Seconds

By aanand – Mon, May 11, 2009 11:00 AM EDT

 Recommend 0  Tweet 3

We used [Apache Hadoop](#) to compete in [Jim Gray's Sort](#) benchmark. Jim's Gray's sort benchmark consists of a set of many related benchmarks, each with their own rules. All of the sort benchmarks measure the time to sort different numbers of 100 byte records. The first 10 bytes of each record is the key and the rest is the value. The **minute sort** must finish end to end in less than a minute. The **Gray sort** must sort more than 100 terabytes and must run for at least an hour. The best times we observed were:

Bytes	Nodes	Maps	Reduces	Replication	Time
500,000,000,000	1406	8000	2600	1	59 seconds
1,000,000,000,000	1460	8000	2700	1	62 seconds
100,000,000,000,000	3452	190,000	10,000	2	173 minutes
1,000,000,000,000,000	3658	80,000	20,000	2	975 minutes

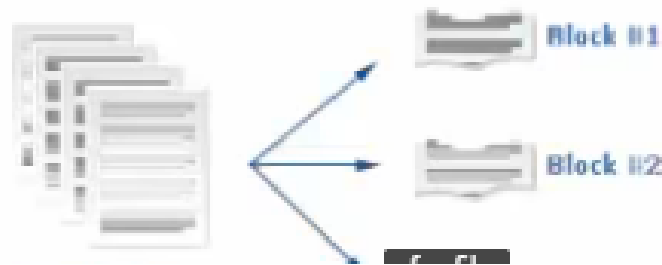
approximately 3000 nodes

Hadoop Architektúra

- Két fő komponens:
 - Distributed File System
 - Hadoop Distributed File System (HDFS)
 - Google File System (GFS)
 - MapReduce Engine / Resource Manager (Yarn)
 - Adat feldolgozó keretrendszer
 - Beépített erőforrás manager és ütemező

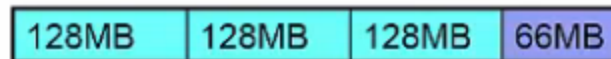
Hadoop Distributed File System (HDFS)

- Meglévő fájlrendszeren működik
 - Nagyobb mennyiségű hiba kezelésére tervezték
 - replikák miatt
- Nagy méretű fájlok tárolására tervezték
 - Folytonos adat olvasására hatékony
 - Véletlen elérés nincs
- A fájlokat blokkokra bontja és azokat tárolja



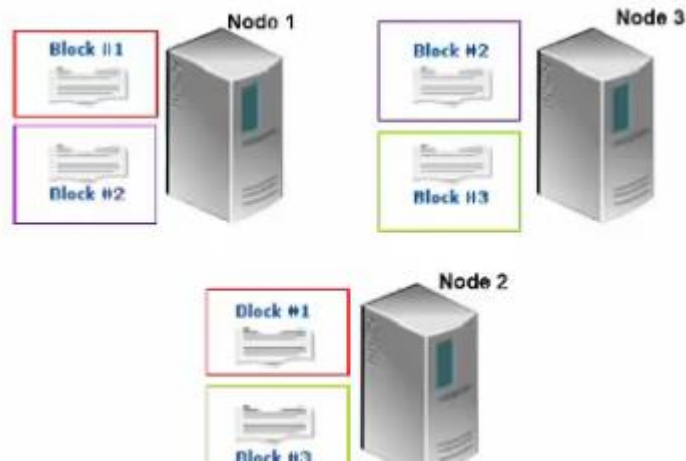
HDFS fájl blokkok

- Nem ugyanaz mint az op. rendszer fájl blokkjai
 - HDFS több operációs rendszeren tárolja a blokkokat
- Blokk méret: 128MB (régábban: 64MB)
- Fájl mérete nagyobb lehet mint bármely lemez kapacitása a klaszterben
 - A fájlok több nodeon tárolódnak blokkként
- Ha kisebb a fájl mint a blokkméret
 - Akkor csak a szükséges részt tárolja



HDFS replika

- A blokkok több node-on vannak tárolva: replikálva
- Ez biztosítja hogy ne legyen adatvesztés hiba esetén
- Default replikaszám: 3
- Szabály: 1 replikának mindig másik rackben kell lennie



Replika készítés folyamata

1. Blokk elküldése egy node-nak
2. A node elküldi a blokkot egy másik rackben lévő node-hoz
3. Ez a node elküldi az adatot nodenak a saját rackjében
4. Minden node visszaigazolja, hogy letárolta az adatot



Klaszter állapotának ellenőrzése

- `hdfs dfsadmin -report` vagy Ambari felület
- Megmutatja a HDFS állapotát
 - Hiányzó blokkok száma
 - Nem replikált blokkok száma

```
C:\Users\Gogo>hdfs dfsadmin -report
A fájlnev, a könyvtárnev vagy a kötetcímké szintaxisa nem megfelelő.
Configured Capacity: 499739779072 (465.42 GB)
Present Capacity: 63237525417 (58.89 GB)
DFS Remaining: 63236362240 (58.89 GB)
DFS Used: 1163177 (1.11 MB)
DFS Used%: 0.00%
Under replicated blocks: 12
Blocks with corrupt replicas: 0
Missing blocks: 0
Missing blocks (with replication factor 1): 0
Pending deletion blocks: 0

-----
Live datanodes (1):

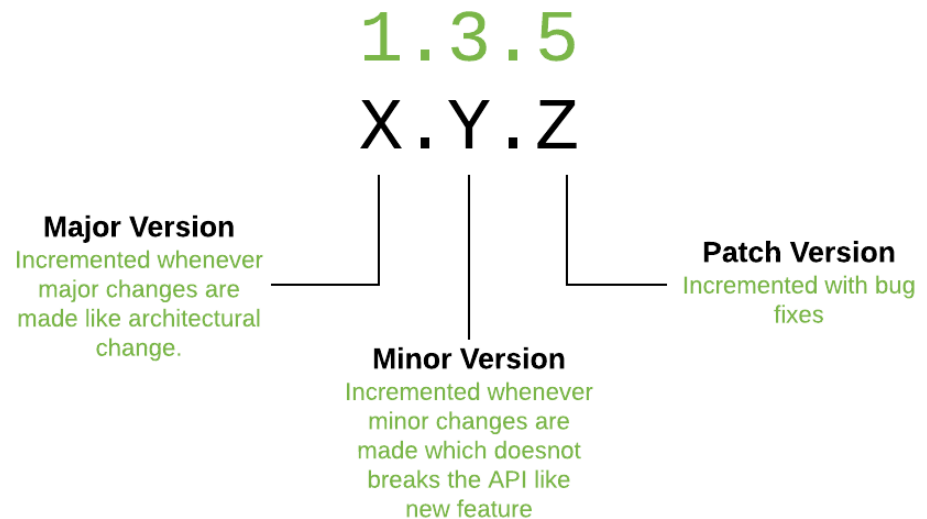
Name: 127.0.0.1:50010 (wit-ams-cloudservice.cloudapp.net)
Hostname: Gogo-PC
Decommission Status : Normal
Configured Capacity: 499739779072 (465.42 GB)
DFS Used: 1163177 (1.11 MB)
Non DFS Used: 436502253655 (406.52 GB)
DFS Remaining: 63236362240 (58.89 GB)
DFS Used%: 0.00%
DFS Remaining%: 12.65%
Configured Cache Capacity: 0 (0 B)
Cache Used: 0 (0 B)
Cache Remaining: 0 (0 B)
Cache Used%: 100.00%
Cache Remaining%: 0.00%
Xceivers: 1
Last contact: Thu Sep 10 13:34:08 CEST 2020
Last Block Report: Thu Sep 10 13:33:27 CEST 2020
```

Verziók

Download

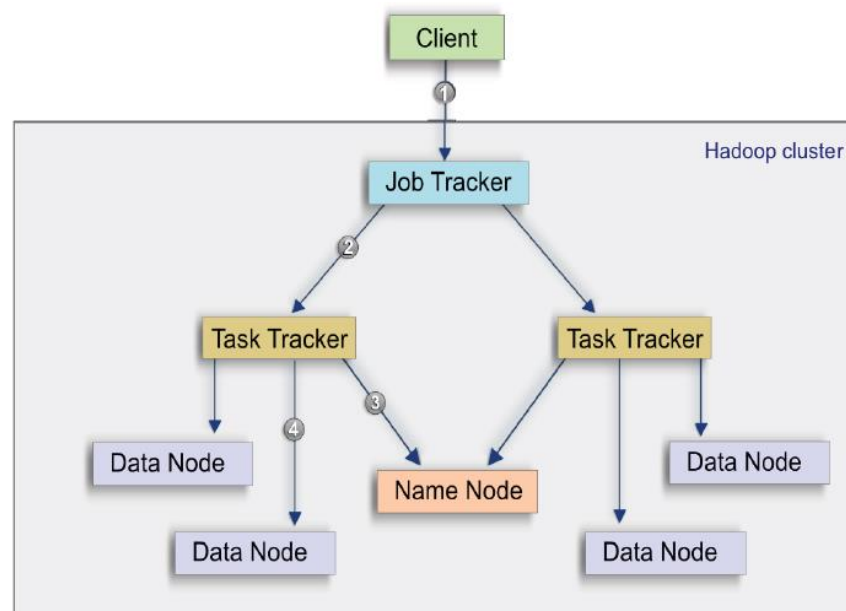
Hadoop is released as source code tart using GPG or SHA-512.

Version	Release date
3.1.4	2020 Aug 3
3.3.0	2020 Jul 14
2.10.0	2019 Oct 29
3.2.1	2019 Sep 22
2.9.2	2018 Nov 19



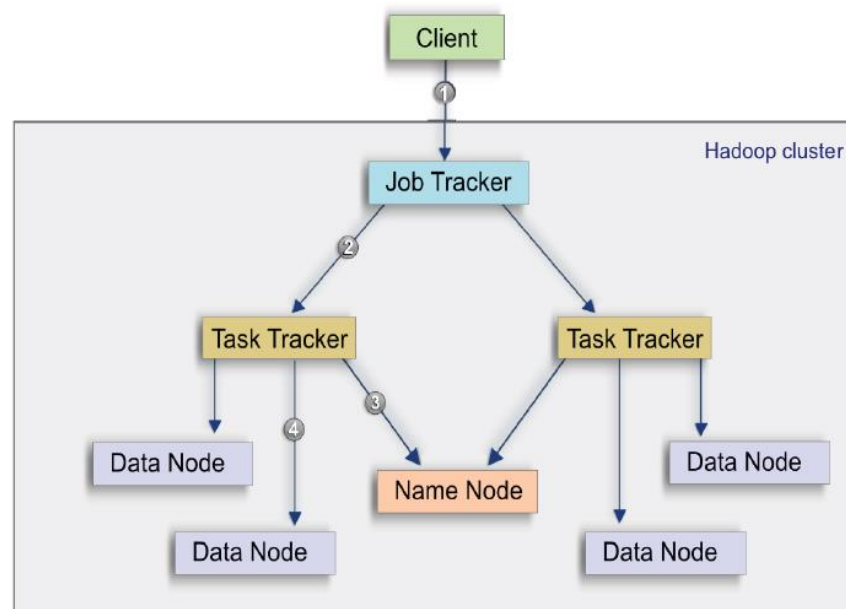
Architektúra (Hadoop 1.x)

- JobTracker:
 - MapReduce Jobok kezelése
 - Klienssel való kommunikáció
 - Feladat kiosztása a TaskTrackereknek
- Task Tracker
 - Adott task folyamat elvégzése:
 - Map
 - Reduce
 - Shuffle



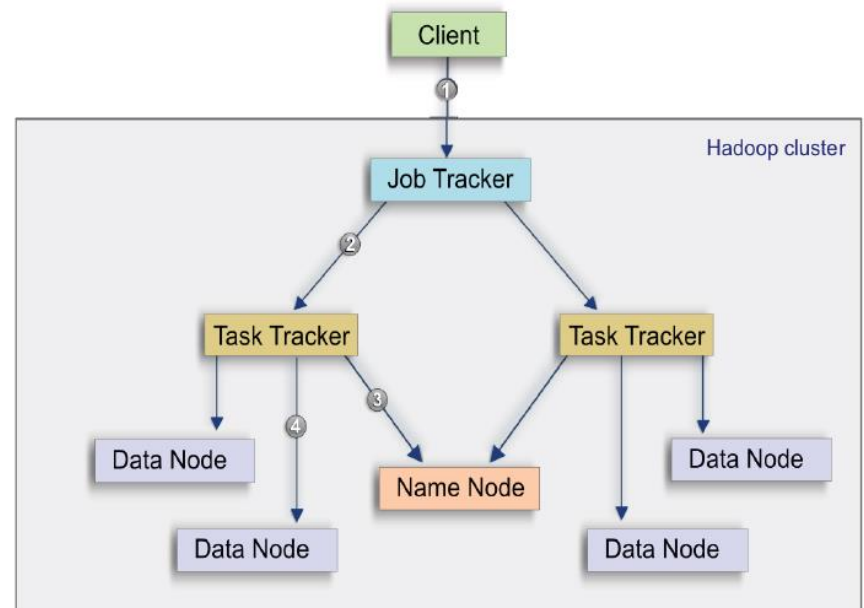
Architektúra (Hadoop 1.x)

- NameNode:
 - HDFS Könyvtár struktúra kezelés
 - Fájlok meglétének ellenőrzése
 - replikaszám
 - hiányzó blokkok
 - Memóriában tartja az információkat
 - Lehet belőle kettő.
- DataNode
 - Tárolja a fájl blokkokat a HDFS-en.



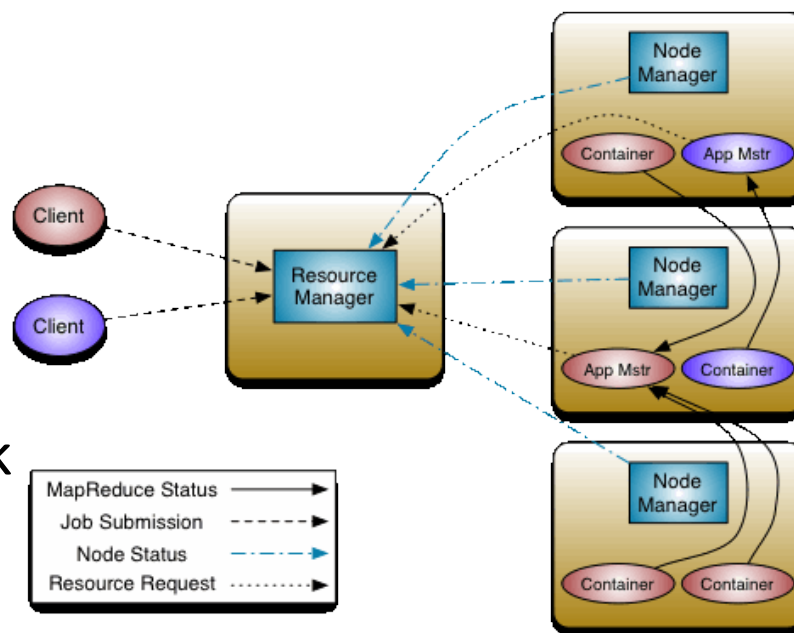
Architektúra (Hadoop 1.x)

- Probléma:
 - Job Tracker szűk keresztmetszet
 - Több kliens kiszolgálásakor túlterhelt
 - Csak MapReduce Jobokat képes kezelni



Architektúra (Hadoop 2.x)

- Resource Manager (YARN)
 - Kapcsolat a kliensekkel
 - Jobok felügyelete
 - A klaszter erőforrásainak ismerete
 - Containerok készítése
- Node Manager
 - Adott node erőforrásainak ismerete
- Application Master
 - Tetszőleges alkalmazás felügyelete és kezelése



Hadoop 3.x



hadoop 3.0 BENEFITS



Support GPUs.



Support
Multiple
Standby
NameNodes



Supports
multiple
NameNodes
for multiple
namespaces.



Storage
overhead
reduced from
200% to 50%.



Intra-Node
Disk
Balancing



MapReduce

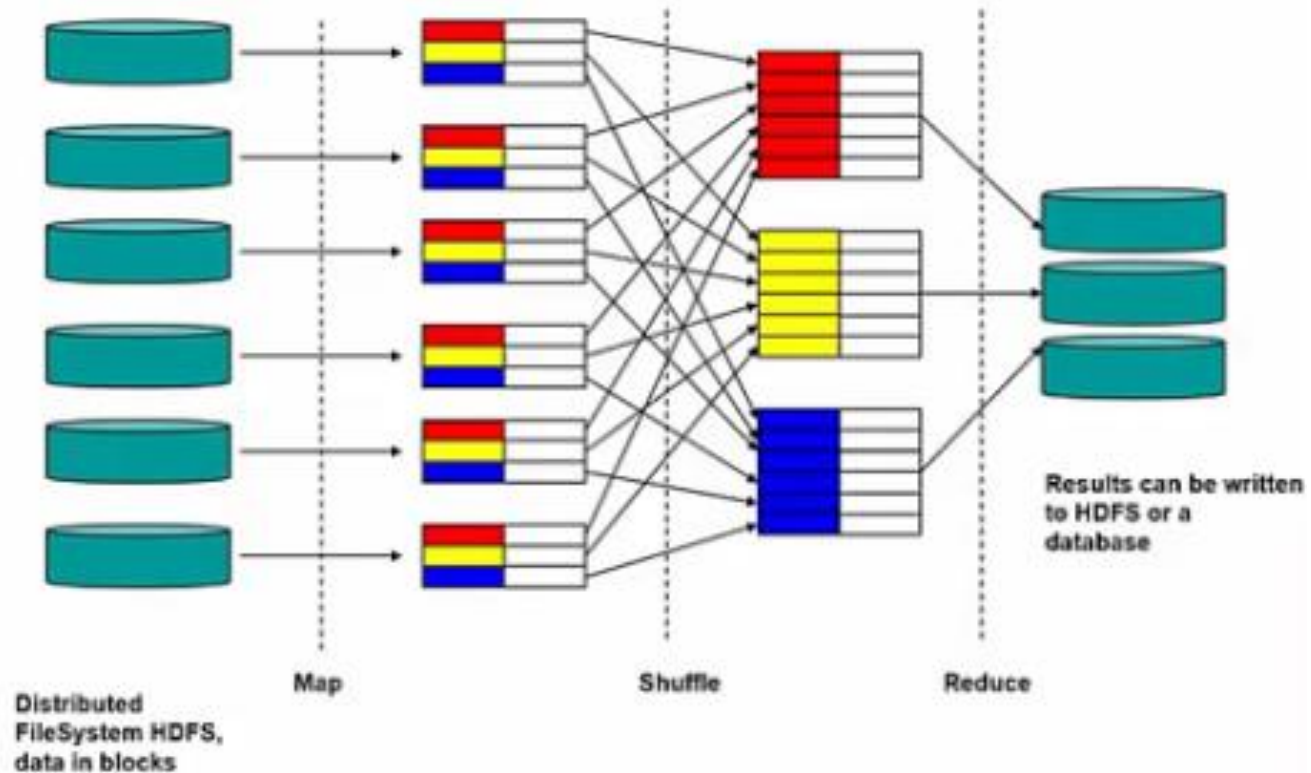
- Nagy mennyiségű adat feldolgozása
- Map
 - Master node szétbontja a problémát kisebb feladatokra
 - Kiosztja ezeket a feladatokat a worker node-oknak
- Reduce
 - Master node összegyűjti a részproblémák eredményét
 - Összegzi a részeredményeket
- A Map és a Reduce folyamatok is képesek párhuzamosan futni

Alap adattípusok

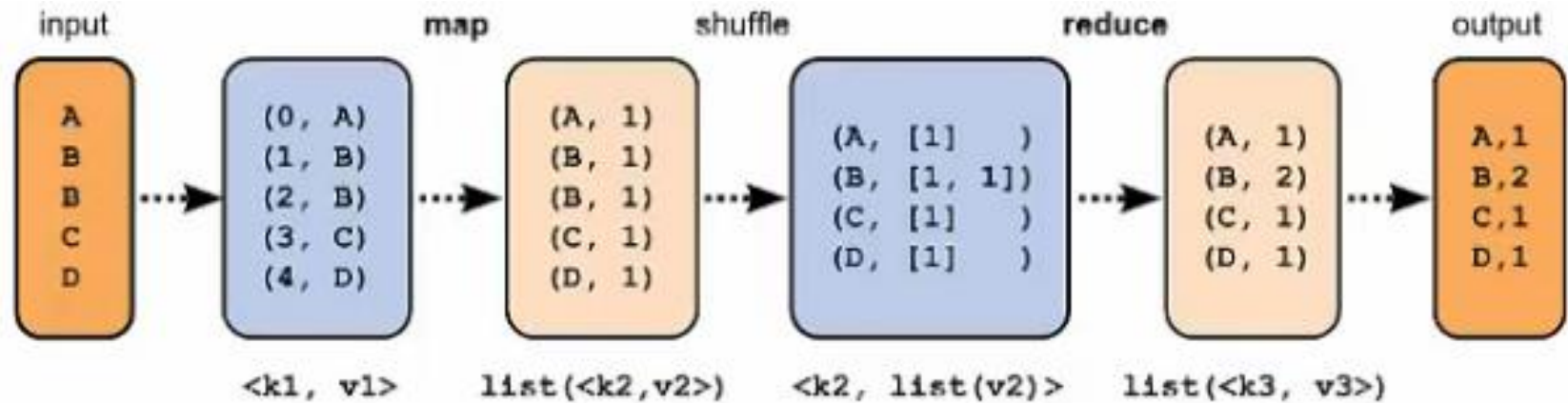
- Kulcs-érték párok
- Listák

	Input	Output
map	<k1, v1>	list(<k2, v2>)
reduce	<k2, list(v2)>	list(<k3, v3>)

MapReduce overview

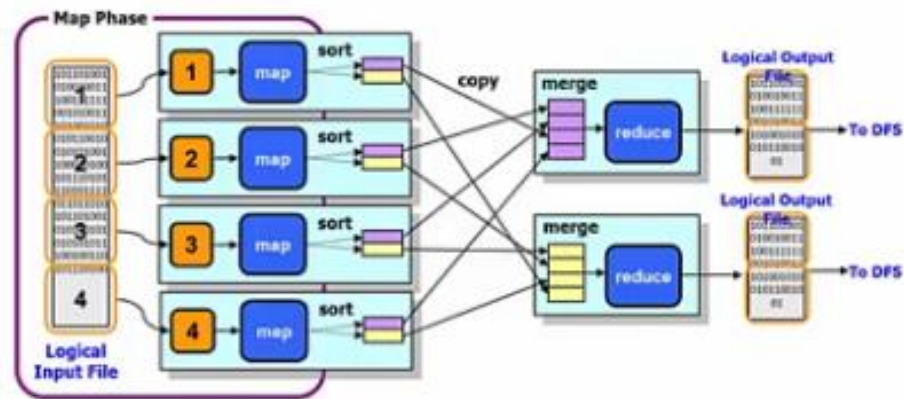


Egyszerű adatfolyam példa



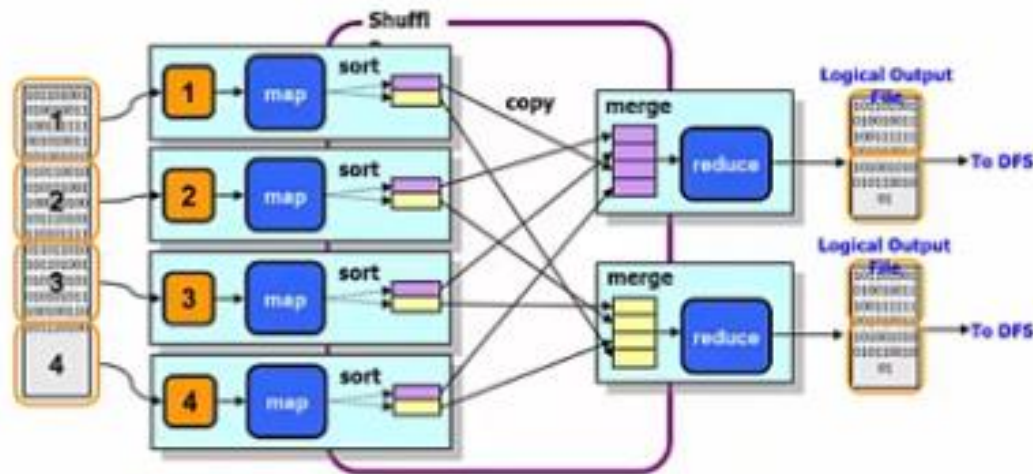
MapReduce – Map fázis

- Mapperek
 - kis program, elosztva a klaszterben, lokális adaton fut
 - az input fájl csak egy kis részén dolgozik
 - minden mapper értelmezi, szűri vagy átalakítja a bemenetet
 - <kulcs,érték> párokat hoz létre



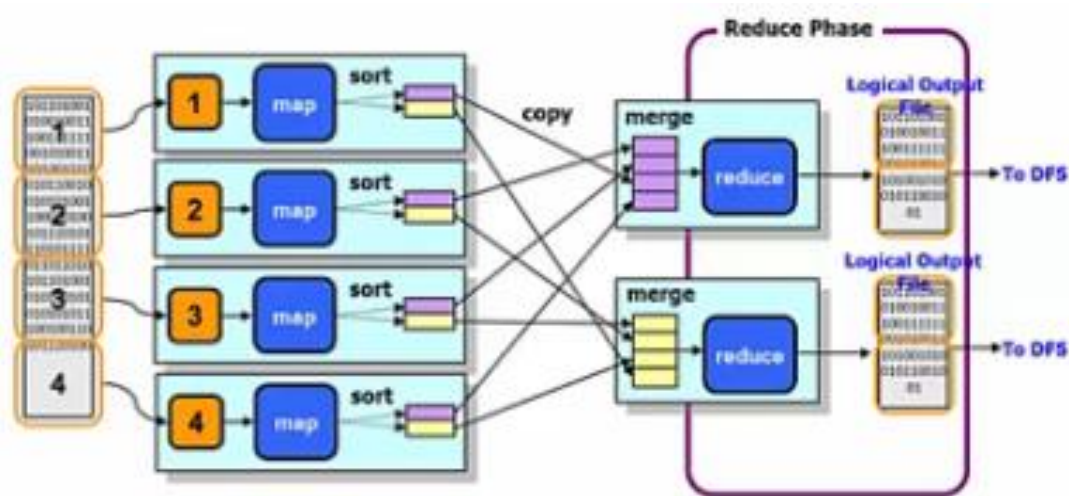
MapReduce – Shuffle fázis

- Minden mapper kimenete lokálisan csoportosítva van kulcs alapján
- Minden egyes kulcshoz kiválasztunk egy node-t
- Minden ilyen adatmozgatást a MapReduce végez



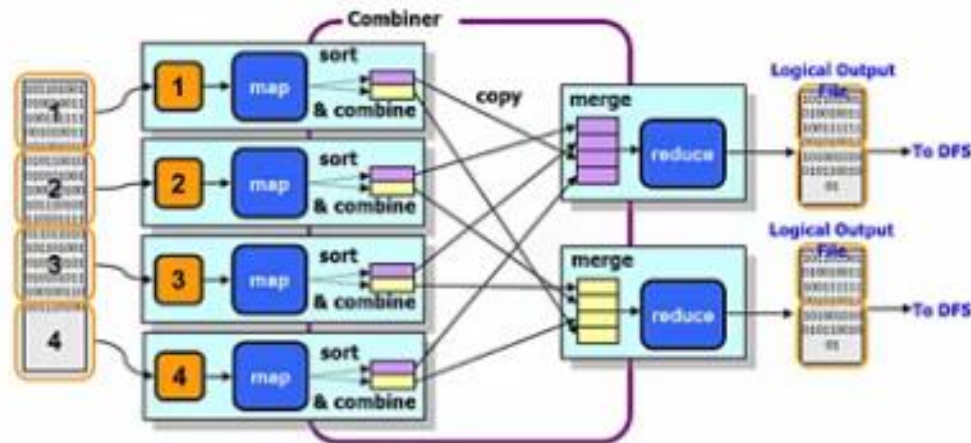
MapReduce – Reduce fázis

- Reducerek
 - kis programok, amelyek aggregálják a hozzájuk tartozó kulcshoz rendelt értékeket
 - Minden reducer saját fájlba írja a kimenetét



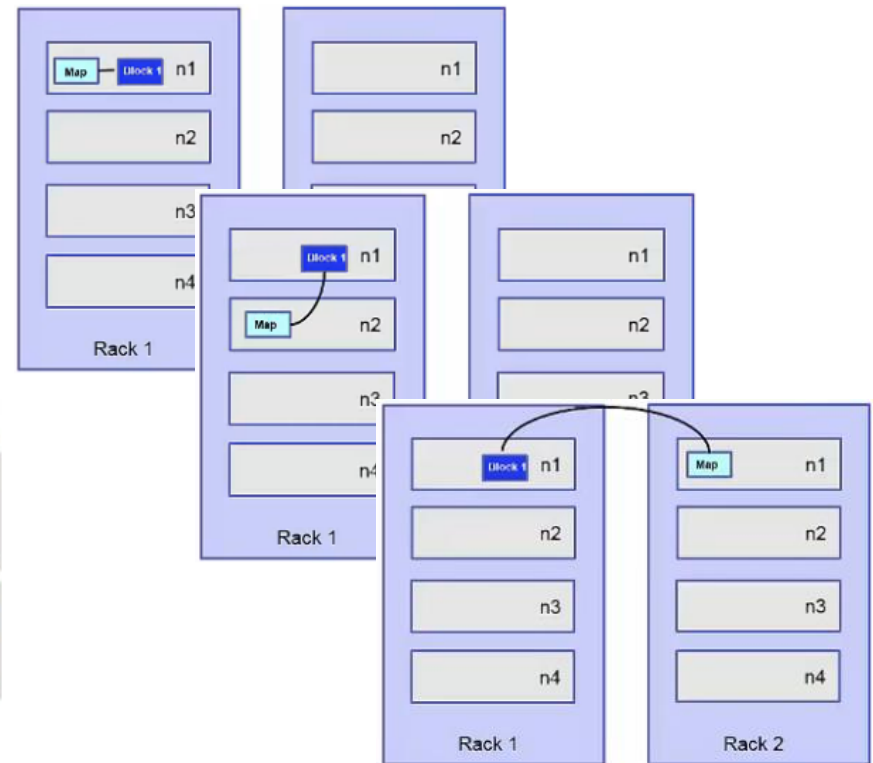
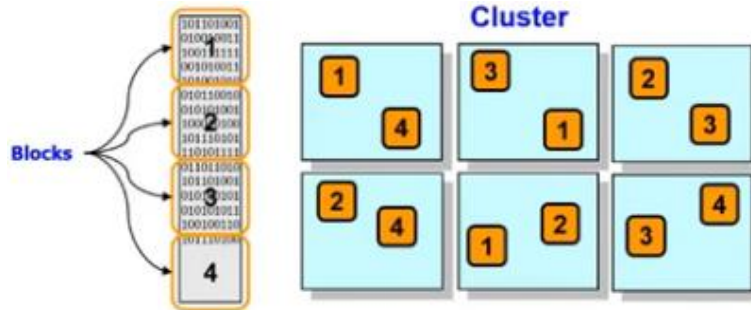
MapReducer – Combiner (opcionális)

- Minden adatot lerendezünk és csoportosítunk mielőtt eljut a reducer node-hoz. Ahhoz hogy csökkentük a hálózati forgalmat, a map oldalon már tudunk végezni egy előzetes redukálást.



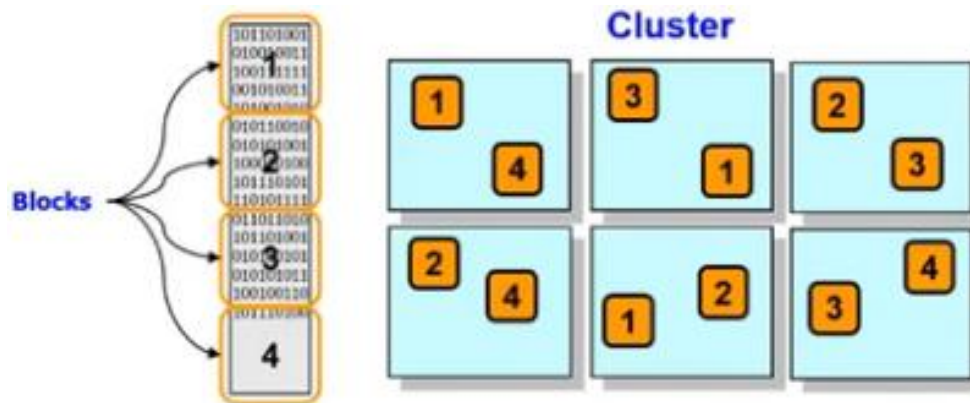
MapReduce – Elosztott fájlrendszer

- Vezérlési szabály
 - az adat a teljes klaszterben van tárolva
 - a programot mozgassuk az adathoz, ne az adatot a programhoz
- Topológia függő futtatás
 1. Ott ahol az adat van
 2. Abban a rackben ahol az adat van
 3. Bármely rackben



Speculative Execute

- Egyes node-ok lassúak lehetnek
- Ugyanaz az adat megtalálható több node-on is.
- Futtassuk ugyanazt a feladatot ezeken a node-okon.
- A leggyorsabb eredményét fogadjuk el a többit elvetjük.



Word Count példa

- Ebben a példában állatnevek vannak
 - MapReduce automatikus darabolja a fájlt sortörések alapján
 - A fájl két részre lett bontva, és két node-on tárolódik
- Meg szeretnénk számolni milyen gyakran fordulnak elő a nagymacska állatnevek
 - SQL-ben így nézne ki:

```
SELECT COUNT(NAME) FROM animals  
WHERE name IN ("Tiger", "Lion", ...)  
GROUP BY name;
```

Node 1

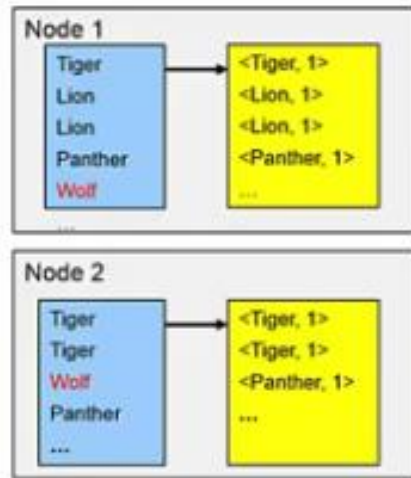
```
Tiger  
Lion  
Lion  
Panther  
Wolf  
...
```

Node 2

```
Tiger  
Tiger  
Wolf  
Panther  
...
```

Map task

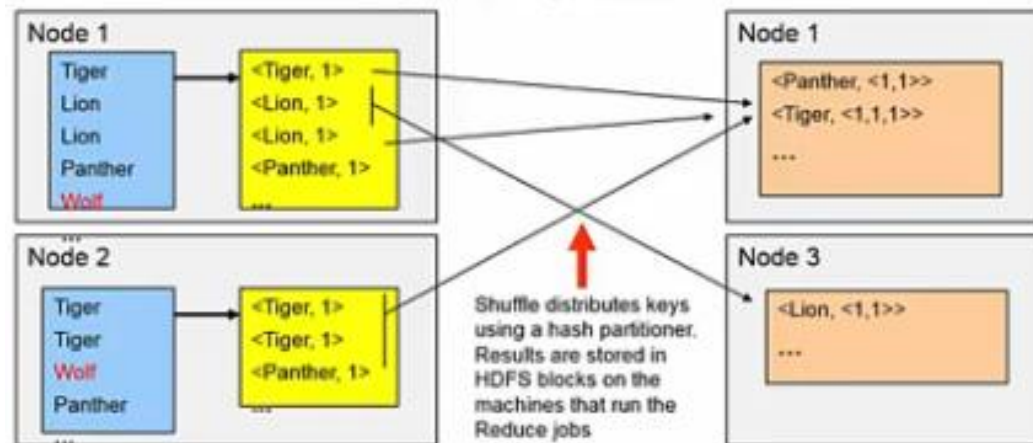
- Két feladtnak lesz a Map folyamatnál:
 - szűrje ki a nem nagy macskákat
 - készítse elő s számoláshoz az adatot:
 - `<Text(name), Integer(1)>`



The Map Tasks
are executed
locally on each
split

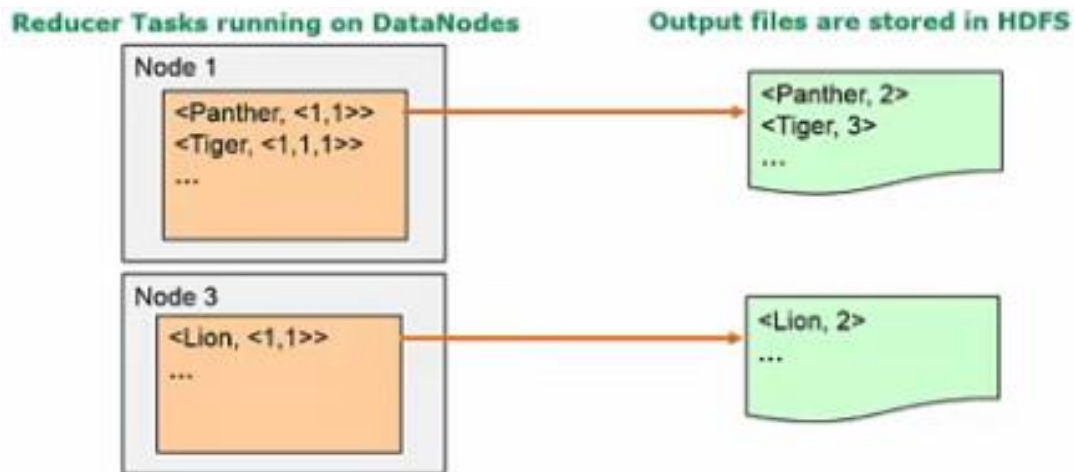
Shuffle

- Shuffle fázis mozgatja az értékeket egy adott kulcs alapján a megfelelő node-okhoz
- Elosztást a Partioner Class végzi (hash elosztás)
- Reducer bármely node-on képes futni (itt node1, node3)
 - Különböző számú Mapper és Reducer task is lehet



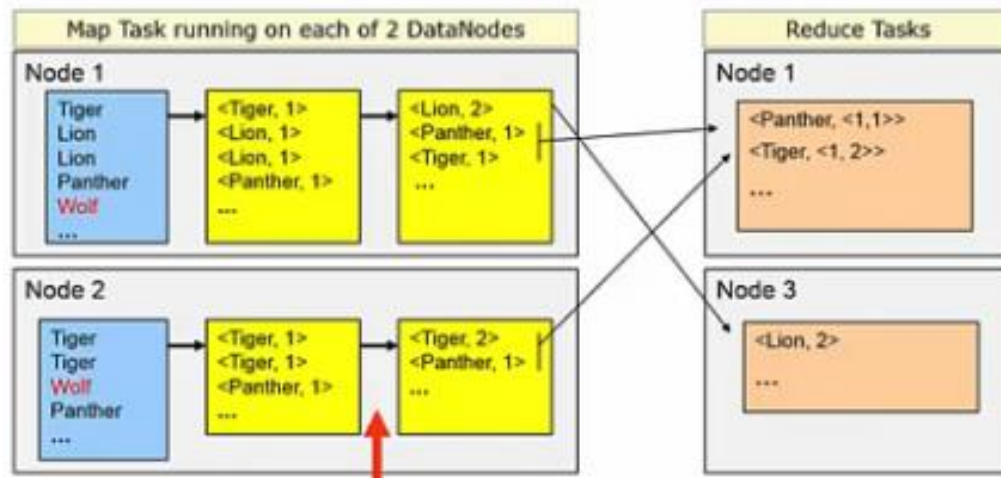
Reduce

- A reducer aggregálja a kulcshoz tartozó értékeket
 - Az eredmény a HDFS-re lesz visszaírva
 - Alapból egy file / reducer
 - Reducer aggregálja az értékeket egy adott kulcshoz (ebben az esetben az állatnevekhez)



Combiner (Optional)

- Combiner
 - opcionális, de gyorsítani tudja a futást
 - aggregációt végez a map folyamat után
 - kevesebb adatot küldünk át a hálózaton a reducereknek
 - Shuffle fázis előtt fut



Writeables

Java Class	Writeable
boolean	BooleanWriteable
byte	ByteWriteable
int	IntWriteable
float	FloatWriteable
long	LongWriteable
String	Text
double	DoubleWriteable

Mapper

```
1 package hu.elte;
2
3 import java.io.IOException;
4
5 import org.apache.hadoop.io.IntWritable;
6 import org.apache.hadoop.io.LongWritable;
7 import org.apache.hadoop.io.Text;
8 import org.apache.hadoop.mapreduce.Mapper;
9
10 public class WCMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
11     private final Text wordKey = new Text("");
12     private final IntWritable wordNum = new IntWritable(1);
13
14     @Override
15     public void map(LongWritable key, Text value, Context context)
16         throws IOException, InterruptedException {
17         String[] fields = value.toString().split(" ");
18         for (String s : fields) {
19             wordKey.set(s);
20             context.write(wordKey, wordNum);
21         }
22     }
23 }
24
```


Reducer

```
1 package hu.elte;
2
3 import java.io.IOException;
4
5 import org.apache.hadoop.io.IntWritable;
6 import org.apache.hadoop.io.Text;
7 import org.apache.hadoop.mapreduce.Reducer;
8
9 public class WCReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
10     private final IntWritable wordNum = new IntWritable(1);
11
12     public void reduce(Text _key, Iterable<IntWritable> values, Context context)
13         throws IOException, InterruptedException {
14         int sum = 0;
15         for (IntWritable val : values) {
16             sum += val.get();
17         }
18         wordNum.set(sum);
19         context.write(_key, wordNum );
20     }
21 }
22
23 }
```

Driver

```
1 package hu.elte;
2
3 import org.apache.hadoop.conf.Configuration;
4 import org.apache.hadoop.fs.Path;
5 import org.apache.hadoop.io.IntWritable;
6 import org.apache.hadoop.io.Text;
7 import org.apache.hadoop.mapreduce.Job;
8 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
9 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
10
11 public class WCDriver {
12
13     public static void main(String[] args) throws Exception {
14         Configuration conf = new Configuration();
15         Job job = Job.getInstance(conf, "WCTest");
16
17         job.setJarByClass(hu.elte.WCDriver.class);
18         job.setMapperClass(hu.elte.WCMapper.class);
19
20         job.setReducerClass(hu.elte.WCReducer.class);
21
22         // TODO: specify output types
23         job.setOutputKeyClass(Text.class);
24         job.setOutputValueClass(IntWritable.class);
25
26         System.out.println(args[0]+" ---- "+args[1]);
27
28         // TODO: specify input and output DIRECTORIES (not files)
29         FileInputFormat.setInputPaths(job, new Path(args[0]));
30         FileOutputFormat.setOutputPath(job, new Path(args[1]));
31
32         if (!job.waitForCompletion(true))
33             return;
34     }
35 }
36
37
```

Köszönöm a Figyelmet!