

A close-up photograph of an open hard disk drive. The image shows the internal platters, the read/write head assembly, and the magnetic recording media. The platters are a light beige color, and the read/write head is positioned above one of them. The overall composition is technical and industrial.

Industrial Machine Learning on Hadoop and Spark

Maks Nakhodnov, Bremen 2025

Distributed computing. Introduction

Plan:

What is big data, and where does it come from?

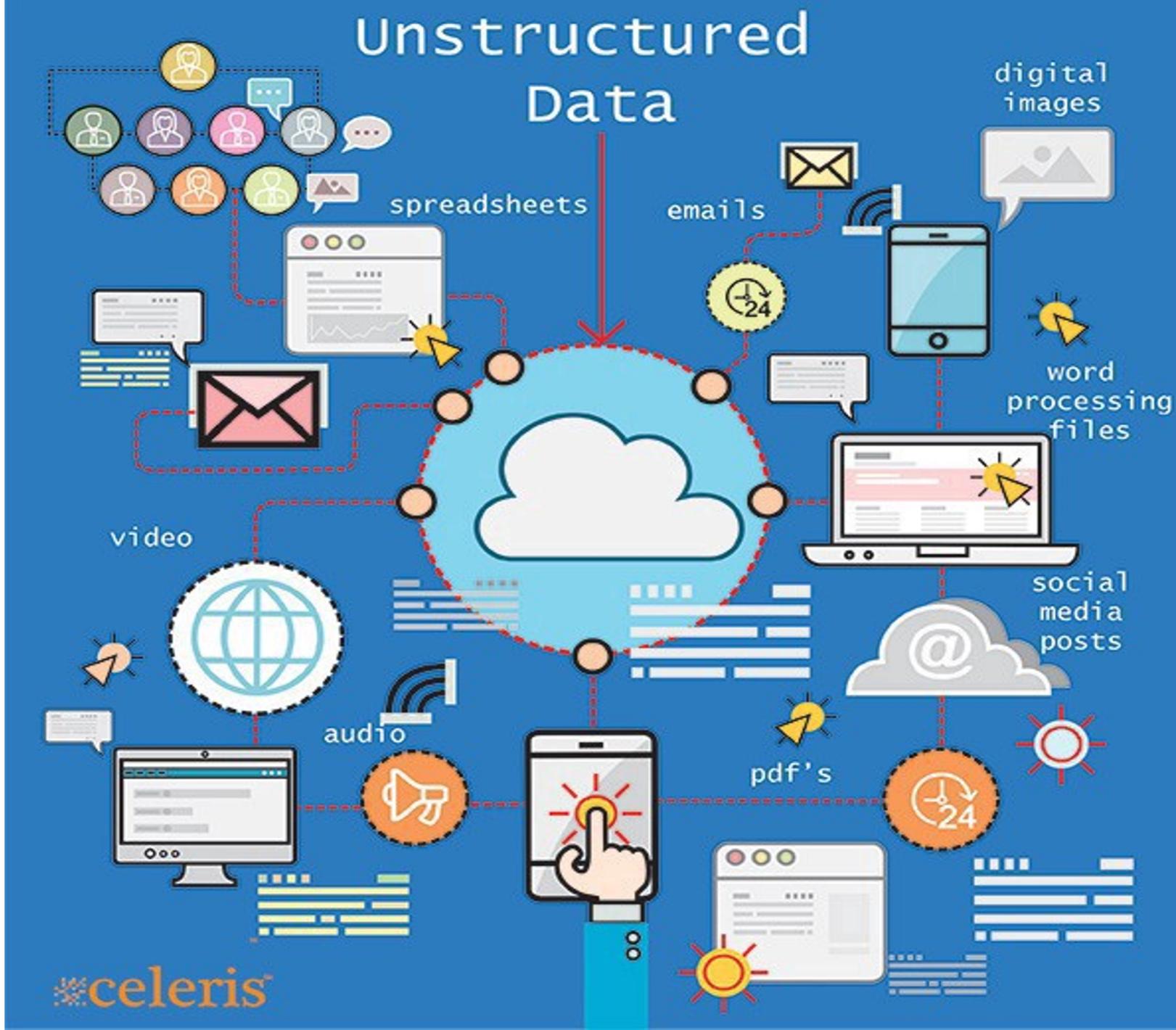
- Industries generating big data. Data explosion.
- Big data – the starting point. 5 main principles.
- How companies handle big data: IaaS / PaaS / SaaS.

Industries generating big data.
Data explosion.

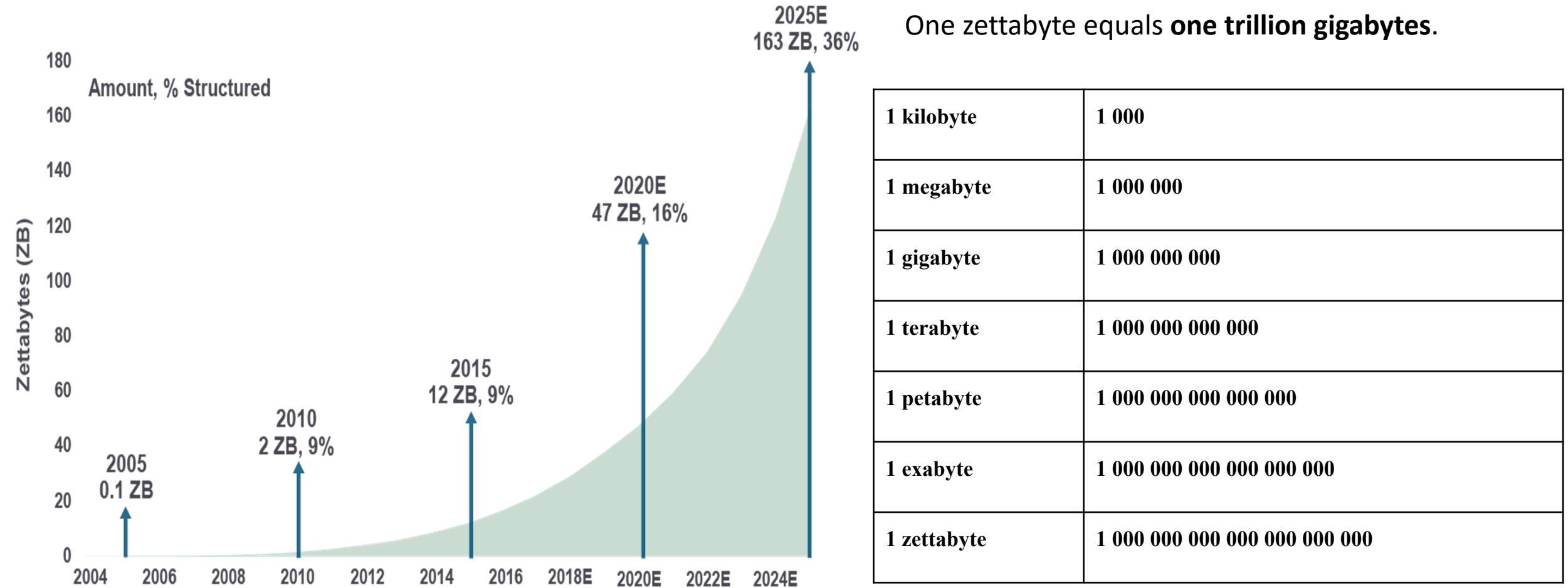
Origins of Big Data

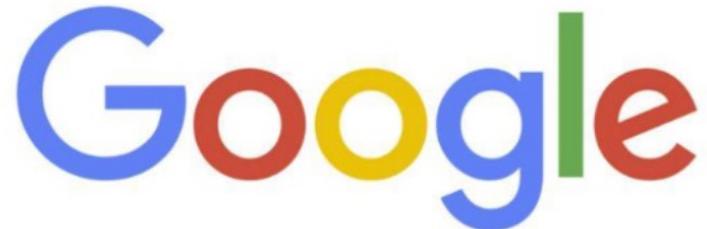
Industries:

- Telecom
- Banking
- Social networks
- Media
- Industry
- Bioinformatics
- Internet of Things (IoT)



Data explosion





Processes 20Pb per day (2008)
Downloads 20B web pages per day (2012)



>10 PB data,
75B DB (6/2012)

>100 Pb user data +
500 TB/day (8/2012)

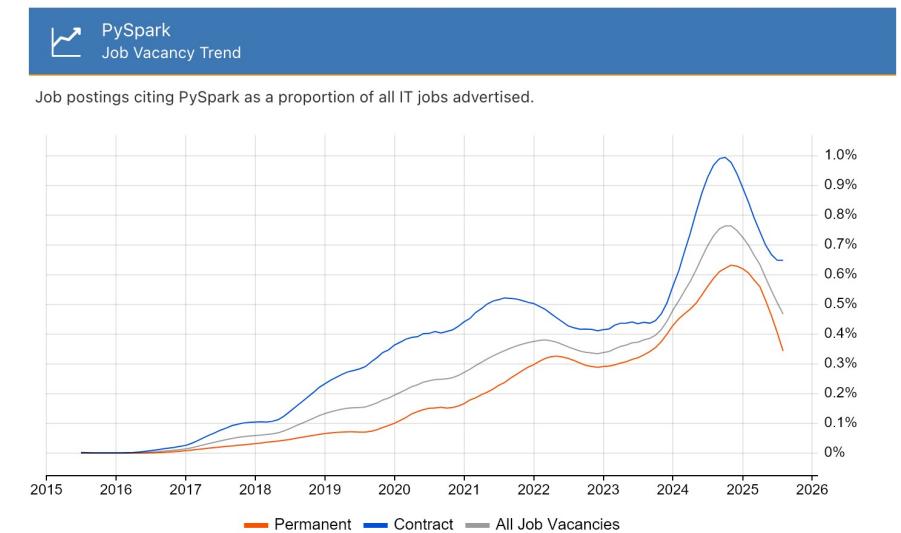
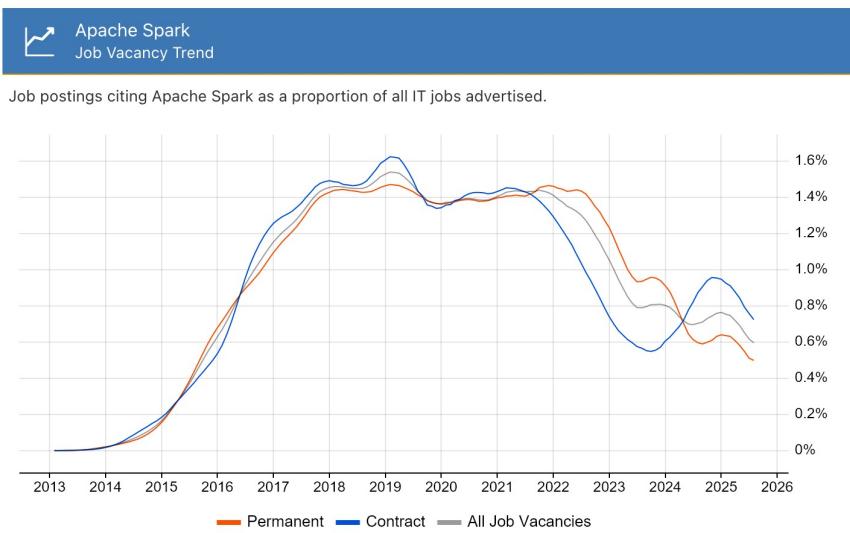
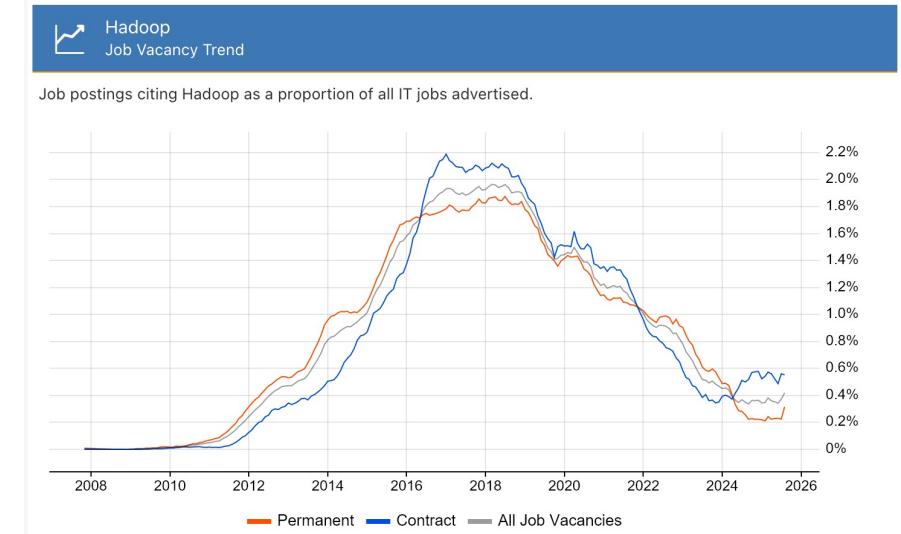
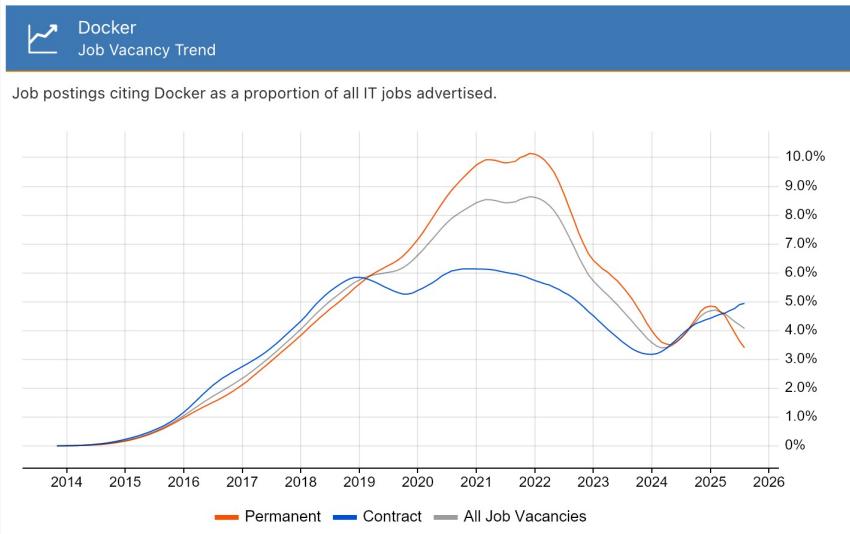


S3: 449B objects (7/2011)

How much is “big data”?



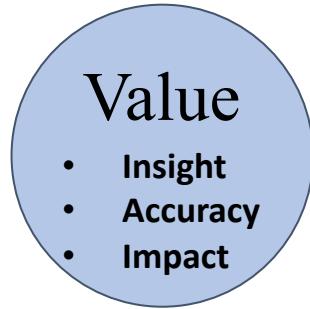
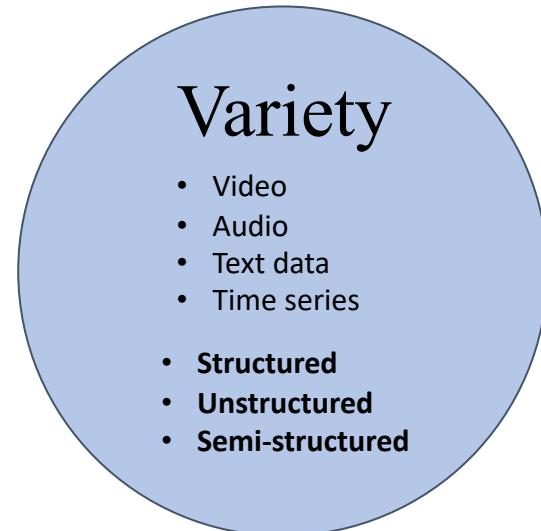
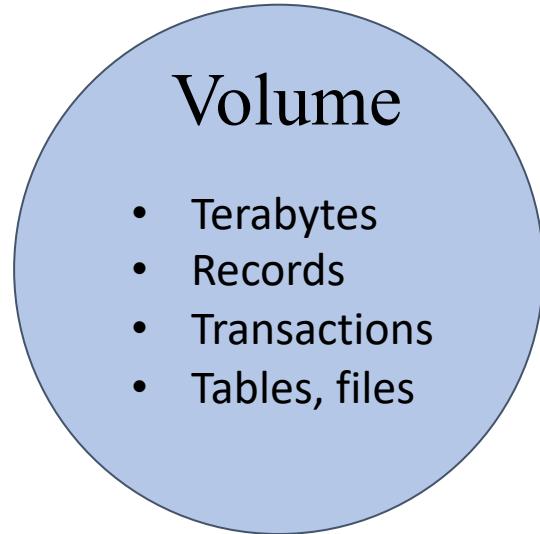
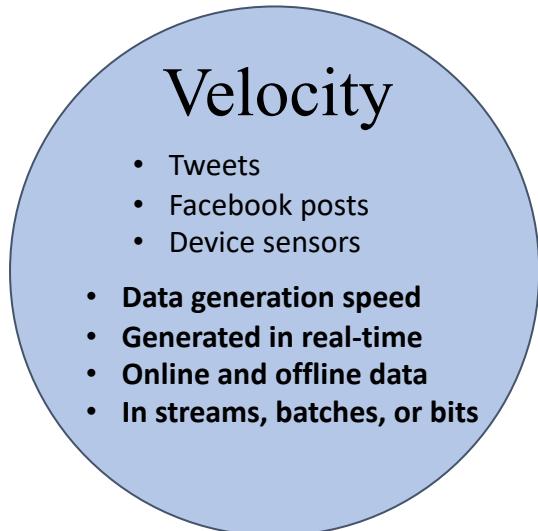
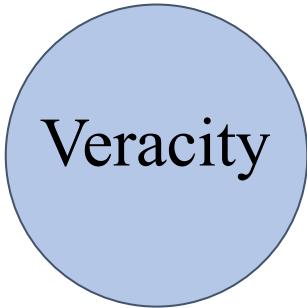
Job trends



Source: <https://www.itjobswatch.co.uk/jobs/uk/>

Big data – the starting point.
5 main principles.

Big data – 5 main principles



How do companies handle big data?

Infrastructure as a Service

Applications

Data

Runtime

Middleware

O/S

Virtualization

Servers

Storage

Networking

Platform as a Service

Applications

Data

Runtime

Middleware

O/S

Virtualization

Servers

Storage

Networking

Software as a Service

Applications

Data

Runtime

Middleware

O/S

Virtualization

Servers

Storage

Networking

You Manage

Other Manages

- **IaaS — Infrastructure as a Service.** Provides infrastructure resources such as virtual servers, storage, and networks ([Google Compute Engine](#), [DigitalOcean](#), [Amazon Web Services \(AWS\)](#), [Cisco Metacloud](#)).
 - Migrating IT systems to the cloud
 - Saving on infrastructure costs
 - Rapid business launch
 - Scaling infrastructure
 - Infrastructure for companies with demand spikes
 - Development and testing
- **PaaS — Platform as a Service.** Provides a platform for building and deploying applications ([Windows Azure](#), [OpenShift](#), [Heroku](#), [Google App Engine](#)).
 - Databases
 - Application development in containers
 - Big data analytics
 - Machine learning
- **SaaS — Software as a Service.** Provides software applications as a service ([Google App Engine](#), [Dropbox](#), [JIRA](#)).
 - Email
 - CRM systems
 - Task planners
 - Website builders

The first key question: how to store big data?

Filesystem

The main functions of a filesystem are:

- Storing and organizing data on a storage device as files
- Determining the maximum supported data volume on the storage device
- Creating, reading, and deleting files
- Assigning and modifying file attributes (size, creation and modification time, owner and creator, read-only, hidden, temporary, archive, executable, maximum filename length, etc.)
- Defining the structure of a file
- Searching for files
- Organizing directories for logical file organization
- Protecting files in case of system failure
- Protecting files from unauthorized access and modification

Linux File System Directories

/bin: Where Linux core commands reside like ls, mv.

/boot: Where boot loader and boot files are located.

/dev: Where all physical drives are mounted like USBs DVDs.

/etc: Contains configurations for the installed packages.

/home: Where every user will have a personal folder to put his folders with his name like /home/likegeeks.

/lib: Where the libraries of the installed packages located since libraries shared among all packages, unlike Windows, you may find duplicates in different folders.

/media: Here are the external devices like DVDs and USB sticks that are mounted, and you can access their files from here.

/mnt: Where you mount other things Network locations and some distros, you may find your mounted USB or DVD.

/opt: Some optional packages are located here and managed by the package manager.

/proc: Because everything on Linux is a file, this folder for processes running on the system, and you can access them and see much info about the current processes.

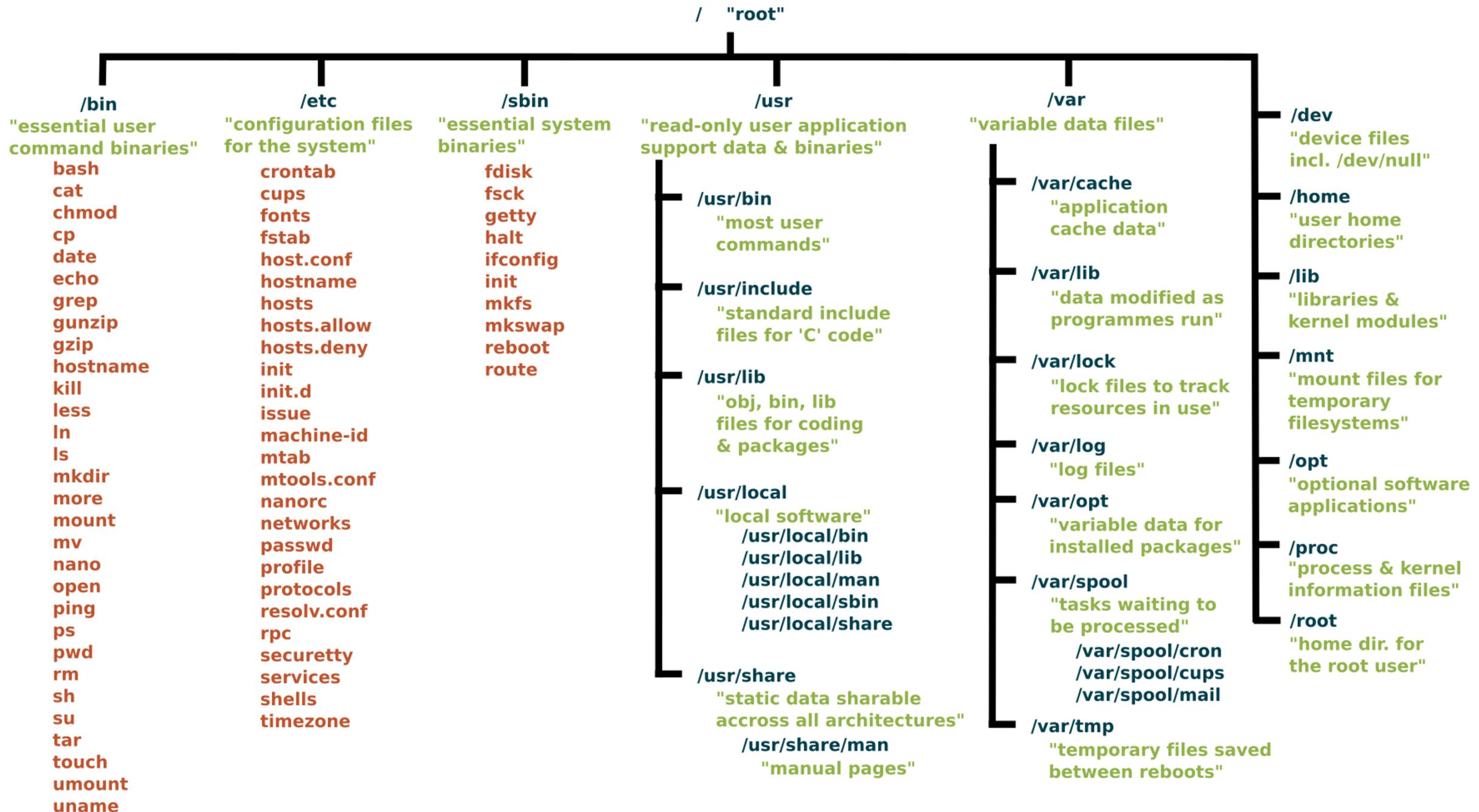
/root: The home folder for the root user.

/sbin: Like /bin, but binaries here are for root user only.

/tmp: Contains the temporary files.

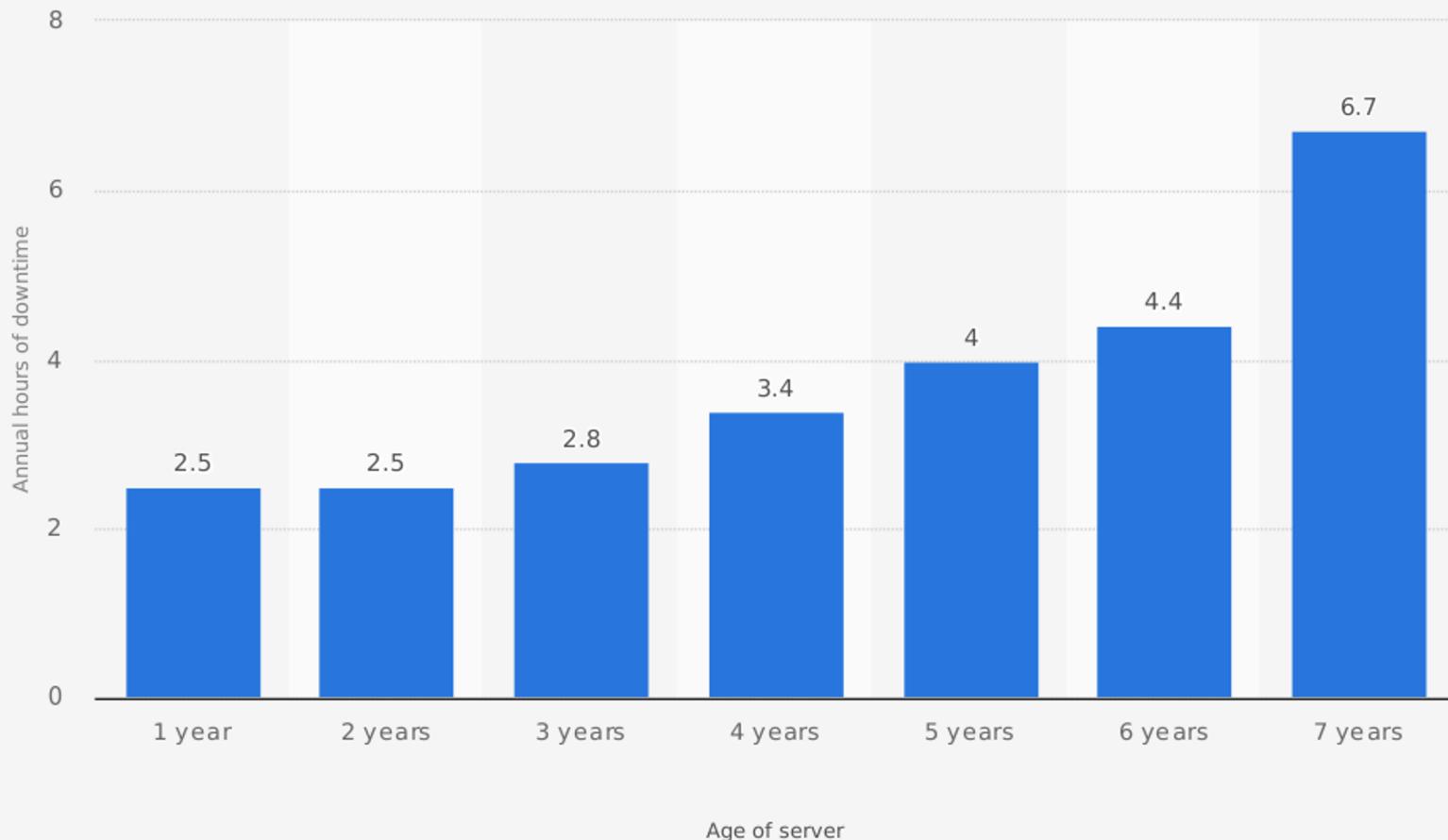
/usr: Where the utilities and files shared between [users on Linux](#).

/var: Contains system logs and other variable data.



Fault tolerance in data storage and computation

Annual number of server downtime hours based on server age, as of 2015



Source
IDC
© Statista 2018

Additional Information:
Worldwide; 2015

The probability that a failure will occur in the next hour:

$$\begin{aligned}P &= 2.5 / (24 * 365) \\&= 0.00028\end{aligned}$$

$$\begin{aligned}P(\text{will not fail}) \\&= (1 - P) = 0.9997\end{aligned}$$

A cluster of 1000 machines

The probability that one of the servers will fail in the next hour:

$$1 - 0.9997^{1000} = 0.25$$