

# Data Management for Big Data

## *Introduction to databases*

**Andrea Brunello**

[andrea.brunello@uniud.it](mailto:andrea.brunello@uniud.it)



→ Different kinds of data

◆ *There's a lot of different stuff out there*

→ The need for databases

◆ *Why can't we simply store files on the disk?*

→ Abstraction levels and data models

◆ *Things that allow us to define and organize data*

→ History of information systems

◆ *When did it all start?*

- Regardless of the domain considered, data are the lifeblood of corporate information systems
- **Data** can be defined as a collection of raw, unorganised and unanalysed material relating to a phenomenon (e.g., a sensor connected to an engine outputting decimal numbers)
- Through a processing and **interpretation** phase, the data become meaningful **information** for the recipient (e.g., the sensor is a thermometer, thus, numbers are interpreted as a temperature in degrees °C)
- Combining information, intuition and experience yields **knowledge**, which can be compared with knowledge already acquired and which allows us to interpret and guide our actions (e.g., based on previous knowledge, the person schedules a maintenance task on the engine)

Interpretation happens also in Natural Language:

- The word «BAT» alone carries little to no information
- Compare «BAT» paired with the words «CAVERN» and «FLY»; and «BAT» paired with «FIELD», «GAME», and «SWING»
- Note that interpretation requires a human (baseball bat still a baseball bat after humanity disappearance?)

Anyway... we interpret also through context!

Collecting data in **organized structures**, and **establishing links between them**, it's what makes them informative

**Information systems** (not to be confused with computer/IT systems) are those (fundamental) systems that support the storage and management of information of an organization for all its purposes

- Traditionally, data considered by information systems used to be very simple
- Things started to change in the **2000s**, following the rise of the Internet and its applications, such as social networks
- We can broadly distinguish between the following kinds of data:
  - ◆ **Structured** data
  - ◆ **Semi-structured** data
  - ◆ **Unstructured** data



Structured data can be seen as tabular data, represented by **rows** and **columns**

- Each row belonging to a same table has a fixed format, think about an Excel file
- For instance, personal data regarding customers
- Easy to store and process, but they convey a limited amount of information

	A	B	C	D	E	F
1	Name	Surname	Birth date	Gender	Address	Phone number
2	John	Doe	12/07/74	M	660 North Gonzales Ave. Canyon City, CA 91387	202-555-0123
3	Mary	Jane	15/08/90	F	7772 Shore St. Zion, IL 60099	202-555-0176
4	Darell	Smart	07/03/83	M	2 Pine Ave. Royal Oak, MI 48067	202-555-0197
5	Jessica	Miller	03/11/95	F	7015 Wilson St. South Portland, ME 04106	202-555-0197
6	Belinda	Barton	05/12/67	F	22 Wakehurst Street Jackson, NJ 08527	202-555-0197
7	Ned	Fitzgerald	26/10/77	M	12 Sage Rd. Hope Mills, NC 28348	202-555-0197

Semi-structured data do not have a fixed format; still, there is some kind of structuring within them

- They make use of **tags** or other markers that allow to organize the content and establish hierarchies within the data
  - For example: XML files
- 

```
1  {
2      "EMPLOYEES": {
3          "SALES": {
4              "648229": {
5                  "NAME" : "Olivia Johnson"
6                  "DOB" : "1989-08-08"
7              },
8              "648666": {
9                  "NAME" : "Frank Mueller"
10                 "DOB" : "1985-05-11"
11                 "MISC" : "On paternal leave from 2019-01-01 until 2020-01-01"
12             }
13         }
14     }
15 }
```

Unstructured data is **not organized** in any predefined manner

- Free text
- Audio and visual materials
- Difficult to store, index, and analyse



Maura

★★★★★ **Soddisfatta**

Recensito in Italia il 30 ottobre 2022

Colore: Bianco | Nome stile: Cassettiera per l'armadio | **Acquisto verificato**

Buon designe, facile da montare, o cassetti sono capienti e robusti hanno il fondo rinforzato con un pannello estraibile di cartone rivestito. Il montaggio lo si può seguire dalle istruzioni illustrate nel libretto, il piano d'appoggio robusto, andrebbe migliorata la rifinitura. Leggero da spostare misure dopo il montaggio altezza 72,5 cm Larghezza 46,80 cm profondità 30 cm



A DBMS contains information on a particular domain

- Collection of **interrelated data** (database)
- Set of **programs** to manage the data

For example, in the university domain:

- **Data**: students, professors, courses, ...
- **Programs**: add a new student, modify the salary of a professor, enrol a student to a course, ...

*DBMSs can be very large,  
and they typically touch all  
aspects of our life!*



In the early days, data-centric applications were built by means of:

- **Files**, e.g., in “csv” format
- **Programs** specifically designed to access such files

First commercial relational DB in the early 80s, after 10 years of work

- So, the need for a DBMS was seen from the beginning

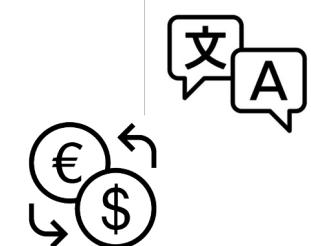
```
Vehicles.txt
Plate,Year,Manufacturer,Model
AX252KF,2015,Ford,Mustang
BK356PQ,2020,Dodge,Charger
FH665HG,2019,Chevrolet,Corvette
CD124LP,2021,Chevrolet,Camaro
```

```
Belongs_to.txt
Plate,Customer_name,Customer_surname
AX252KF,John,Doe
BK356PQ,Mary,Jane
FH665HG,Jessica,Miller
CD124LP,Ned,Fitzgerald
```

```
Customers.txt
Name,Surname,Birth_date,Gender
John,Doe,1974-07-12,M
Mary,Jane,1990-08-15,F
Darell,Smart,1983-03-07,M
Jessica,Miller,1995-11-03,F
Belinda,Barton,1967-12-05,F
Ned,Fitzgerald,1977-10-26,M
```

## Data redundancy and heterogeneity

- Multiple file formats and programming languages
- Duplication of information in different files
  - ◆ Possible data inconsistency!



## Difficulty in accessing data

- Need to write a new program to carry out each new task
- Low interactivity



## Constraints management

- Domain constraints (e.g., car value  $\geq 0$ ) become «buried» in program code rather than being stated explicitly
- Hard to add new constraints and manage existing ones

## Atomicity of updates

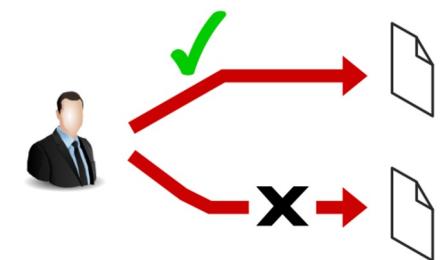
- Failures may leave the system in an inconsistent state
- E.g., a transfer of funds from one account to another should be either carried out in its entirety or not happen at all

## Concurrent access by multiple users

- What happens if two users try to access and modify the same data?
- E.g., two users booking the last airline ticket

## Security problems

- Grant users access to some, but not all, data



*(Relational) DBMSs provide solutions  
to all these problems!*

→ **Large** amounts of data

→ **Global** data

- ◆ They are relevant for a vast array of users and applications
- ◆ The database is typically not tied to a particular user or application, as standard programs are

→ **Persistent** data

- ◆ The database is there independently from the interacting users or applications
- ◆ Data lives on its own unlike, for instance, program variables

- To be capable of dealing with all the previous stuff, a DBMS must be a very complex piece of software
- Thus, to allow for its development and management, it is typically organized following a **stratified approach**
- This means that the DBMS functionalities are arranged into **different hierarchical levels**
  - ◆ each level encapsulates/hides the details it deals with,
  - ◆ it offers its functionalities to the higher levels,
  - ◆ and it can rely on the functionalities made available by the lower levels, without caring about their implementation (**abstraction**)

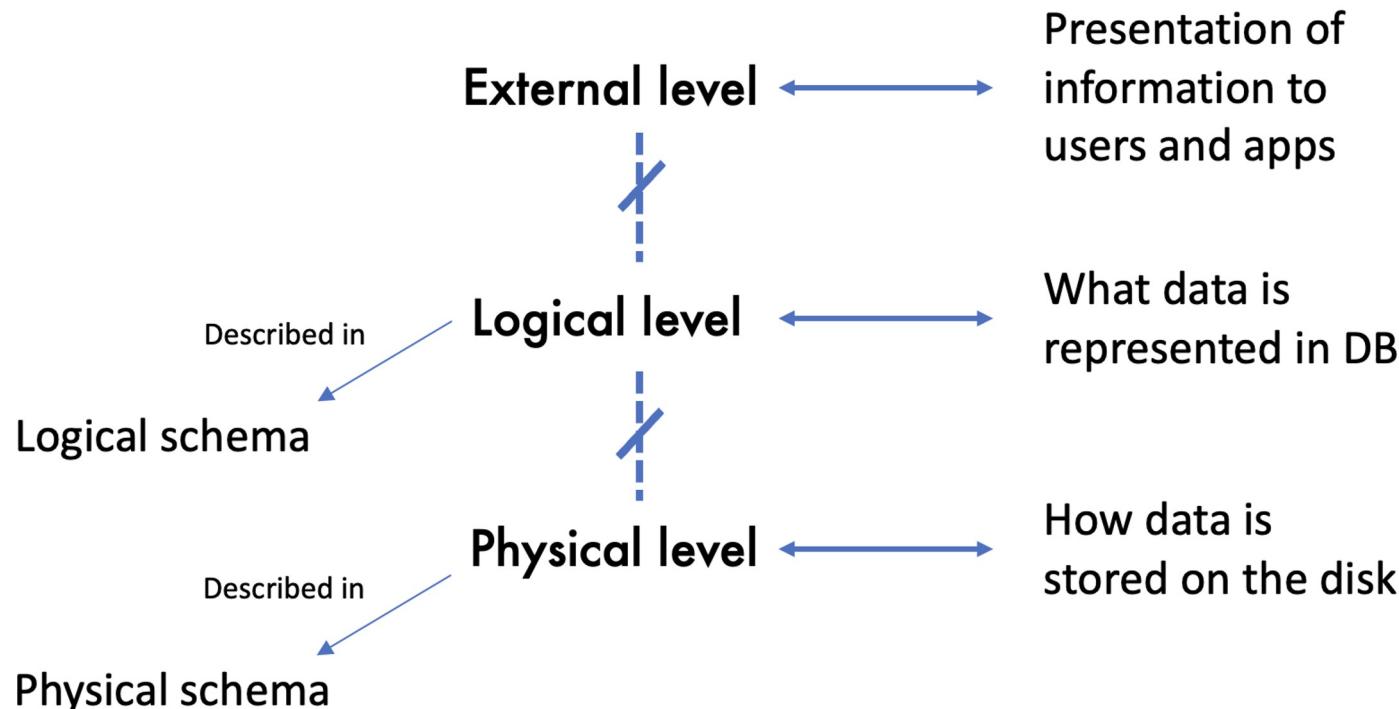
From lowest to highest:

- **Physical level:** *how* data is stored on the disk (data structures, file arrangement, ...)
- **Logical level:** *what* data is stored and represented in the database and their relationships
  - ◆ describes an entire database in terms of a small number of relatively simple structures
  - ◆ e.g., data is stored within tables, or graphs
  - ◆ they are *tables/graphs*, we do not care about how they are stored on the disk (the physical level deals with it)
- **View/External level:** it manages the presentation of information to users and programs. Views can also hide information to some kinds of users (e.g., show only some parts of a table)

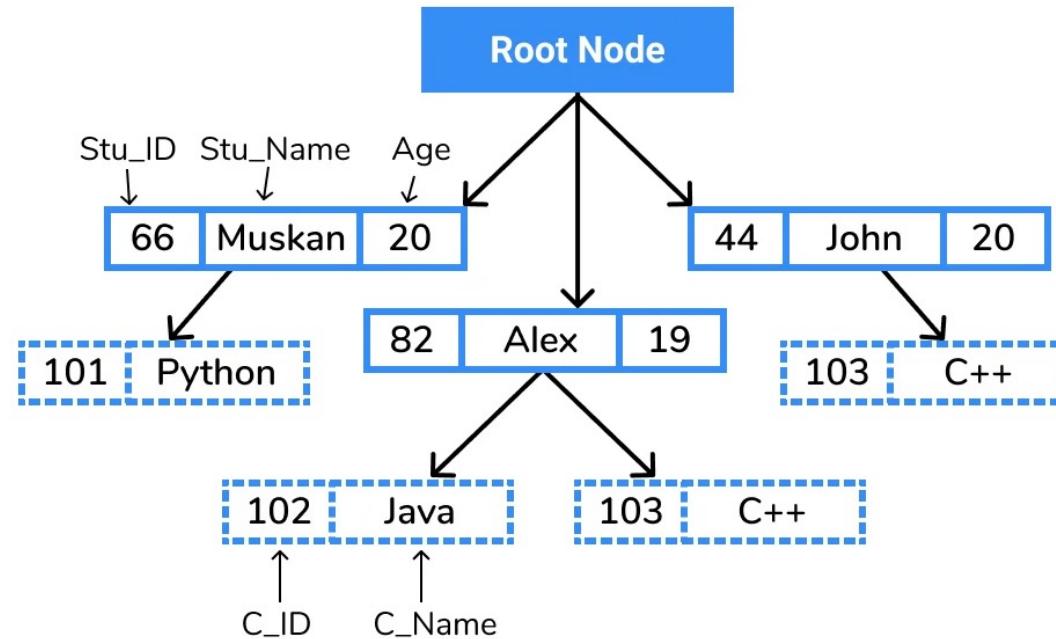
The levels of abstraction of a specific database are described by means of *schemas*:

- **Physical schema:** it lays out how data are stored physically on a storage system in terms of files and data structures
- **Logical schema:** it encodes the overall logical structure of the database
  - ◆ E.g., the database consists of a set of tables that are designed to store information about a set of customers (each characterized by name, address, and phone number) and their accounts in a bank; the relationships among customers and accounts; and, some constraints over them

**Physical data independence:** ability to make changes to the physical level without having to modify the logical one



- To describe the logical schema, a small number of relatively simple structures is used (data model)
- A **data model** is a set of tools (like an alphabet/language) for describing
  - ◆ Data (and its structure)
  - ◆ Data relationships
  - ◆ Data constraints
- There exist several logical data models, including
  - ◆ The **relational data model**, which considers rows and tables (records and relations)
  - ◆ Graph data model, that deals with nodes and vertices
  - ◆ Object-based data models
  - ◆ Hierarchical model

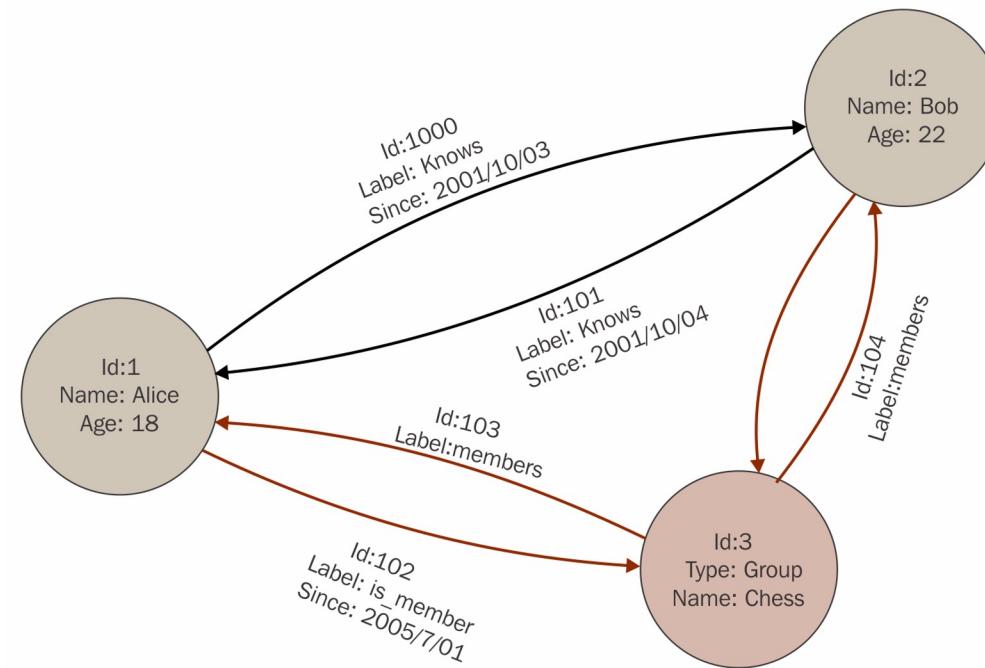


Records physically point to other records; only «parent-child» relationships are allowed

Evolution: network model structure (*The Programmer as a Navigator*, Bachman)

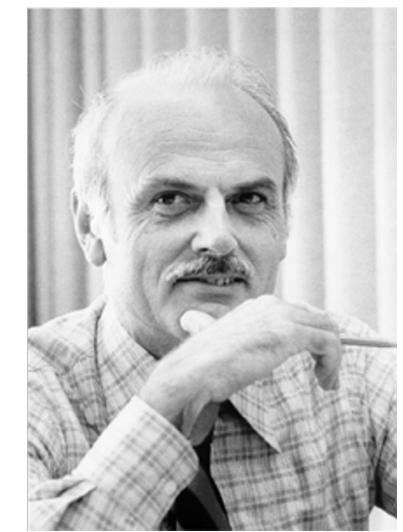
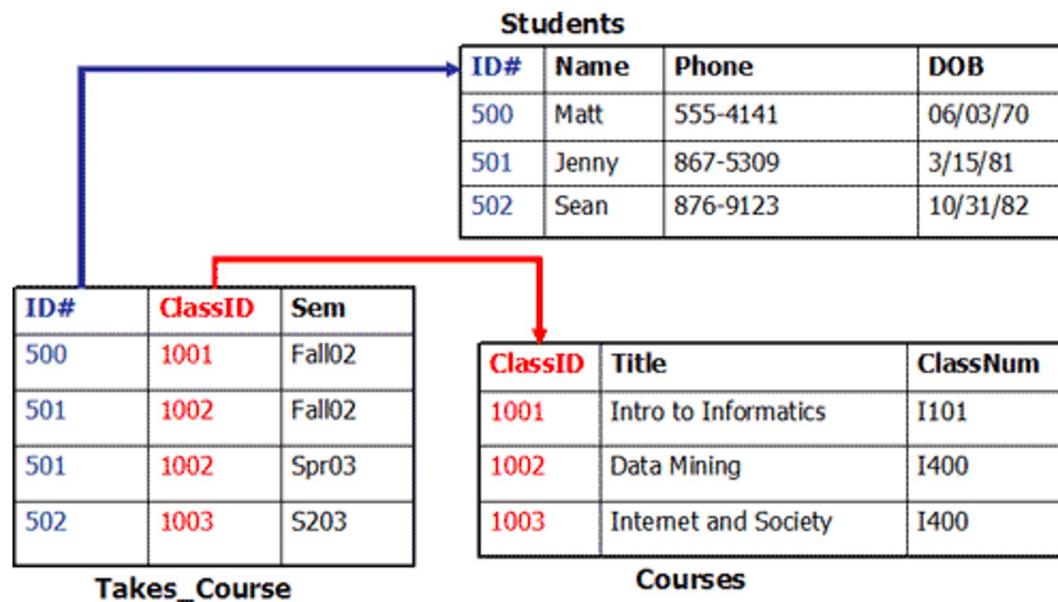


Nodes and arcs to represent data and relationship among the data



Records *physically* point to other records; more flexible structure for the nodes and the relationships

Tables to represent data and relationship among the data



The two data models are equivalent pertaining to their data representation capabilities

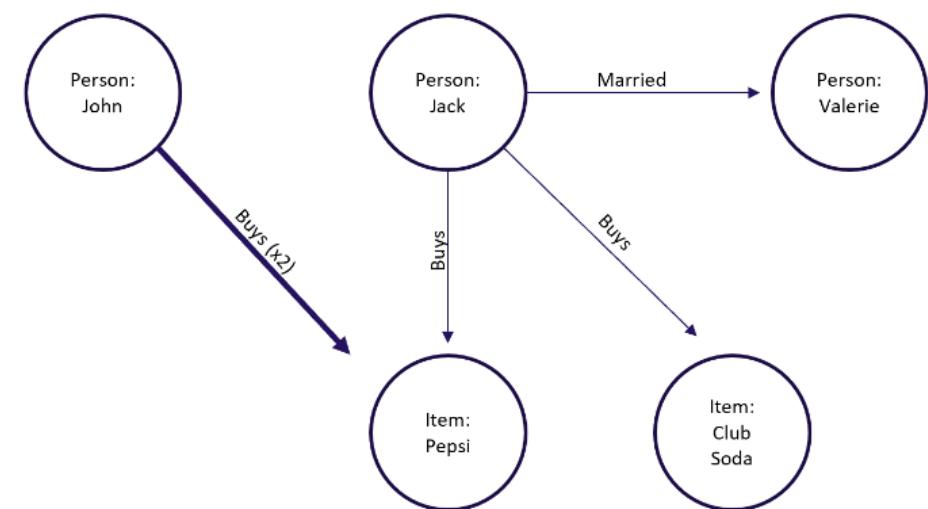
Sales		
Customer	Item	Time
0001	1A	20:34
0001	1A	21:15
0003	2A	21:16
0002	1A	21:16
0002	5C	

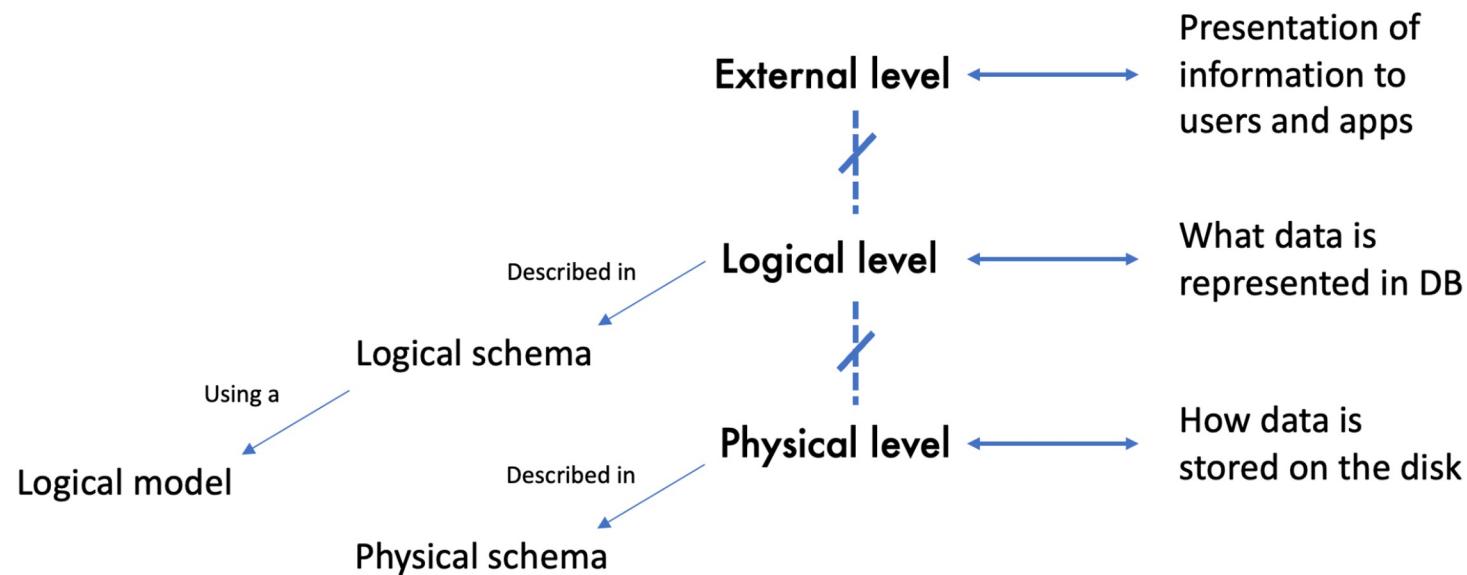
  

Customer	
Name	CustID
John	0001
Jack	0002
Ted	0003
Ken	0004
Valerie	0005

Inventory	
Description	SKU
Pepsi	1A
Club Soda	2A
.	.
.	.
Diet Coke	5C



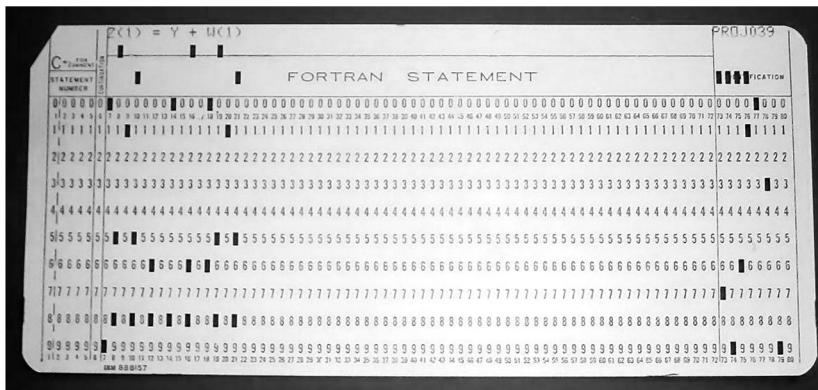


- Every kind of organization has an information system which organizes and manages data to support relevant business goals
- Information systems pre-date computers and even electricity, and are not necessarily automated systems
- For instance, think about accounting and bookkeeping, or the civil registry

The image shows an open ledger with two pages of handwritten entries. The left page is for "D<sup>r</sup> Vincent" and the right page is for "Stillwill C<sup>r</sup>". Both pages show dates, descriptions of transactions, and amounts in pounds, shillings, and pence. The ledger includes various entries such as "By Sandries", "Mending 10 Pairs Stock", "Taping 2 Days Draying Horses", and "Mending 100 Pairs Stock". The right page also contains a note about settling an account with "Durham" on May 27<sup>th</sup>, 1817.

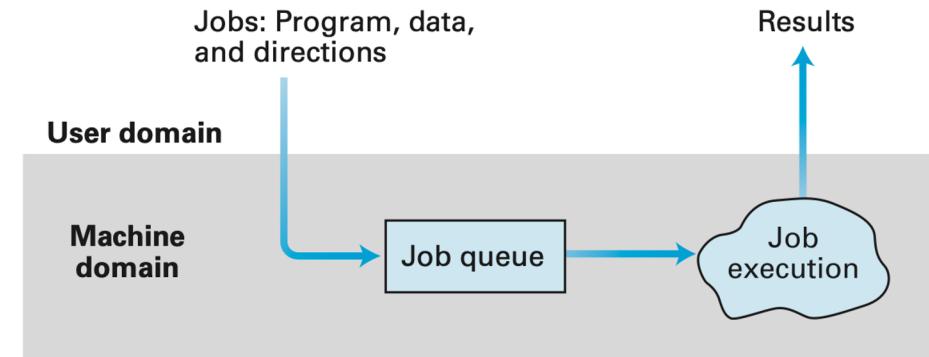
Date	Description	Amount
Jan 28	To Sandries	10 16
30	" Making 10 Pairs Hoses "	5 "
March 19	" 2 Days Draying Horses "	10 "
23	" Taping 8 Pairs Hoses "	2 "
28	" Mending 100 Pairs Stock "	1 "
April 5	" Taping 8 Wall Bars "	2 "
26	" Mending 2 1/2 Pairs .. "	3 "
May 9	" Mending 100 Pairs 1 hours "	2 6
		£ 11 19
1817		
May 27	To Balance Due	1 8
28	" Making 10 Pairs Stock "	6 "
June 12	" Making 10 Pairs Stock "	5 6
21	" Making 8 Pairs "	1 6
August 9	" Taping 10 Pairs Hoses "	2 6
Sept 1	" Rake "	2 "
4	" In Bushel Unions "	1 11
	" 1 Day harrying ) "	9 "
	with Crop )	
	" 1 Day Curving Barn ) "	8 "
	with Crop )	
	" Making 100 Pairs Stock "	5 6
20	" Making 10 Pairs Hoses "	3 6
	" Rake Unions Barnyard "	1 "
24	" Making 100 Pairs Stock "	6 "
Decemb 13	" Mending 100 Pairs Stock "	2 1/2
1818		
Jan 10	" Mending 100 Pairs "	2 1/2
		£ 11 19
D <sup>r</sup> Vincent		
1817		
July 28	By Sandries	9 16 1
May 17	" Porter & Plover	2 5 8
	" Balance Due on account	2 11 1 9
	This day recd & settled all accts	
	& find due two dollars fifty Eight	
	Cnts Durham May 27 <sup>th</sup> 1817	
	Cnts Watrous	
Stillwill C <sup>r</sup>		
1817		
July 1	By Purchased Plaster	4 6
	" Hoof to Stone Ann.	2
	" Hammer	2 "
Oct 1	" Horses & Man	6 "
	" Black Paint	1 10
	" Harness Plastered	3 9
Dec 1	" 8 Bus. Buck wheat	1 19
		3 2 3
Vincent Stillwill		
1818		
July 28	To Sandries	4 9 5
	" Mending 100 Pairs Stock "	1 6
March 3	" Mending Boxes "	2 3
12	" Taping 1 Box "	1 9
20	" Making 100 Pairs Stock "	6 "
		£ 11 19 10

- 1950s and early 1960s
  - Data processing done using magnetic tapes for storage, which provided only **sequential access**
  - Punched cards for input
  - Many dedicated hardware as many formats available
  - Support for tasks such as payroll management

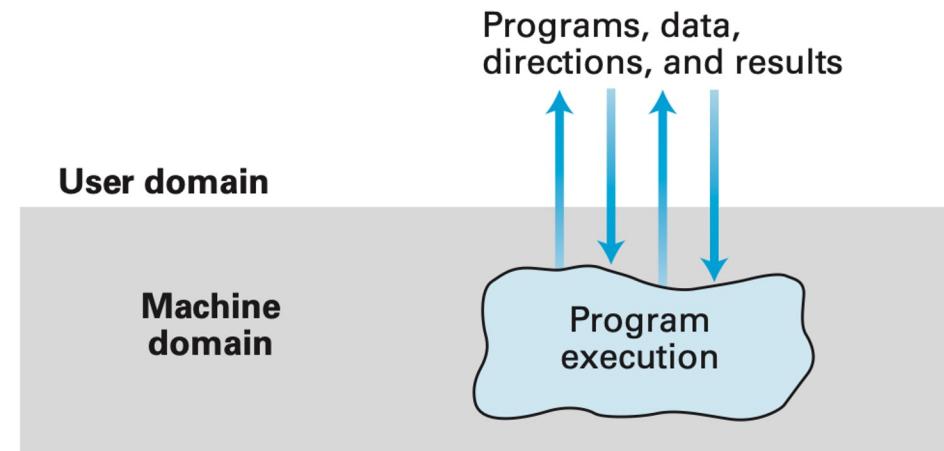


These were the golden years of «batch processing»...

## Batch processing



## Interactive processing



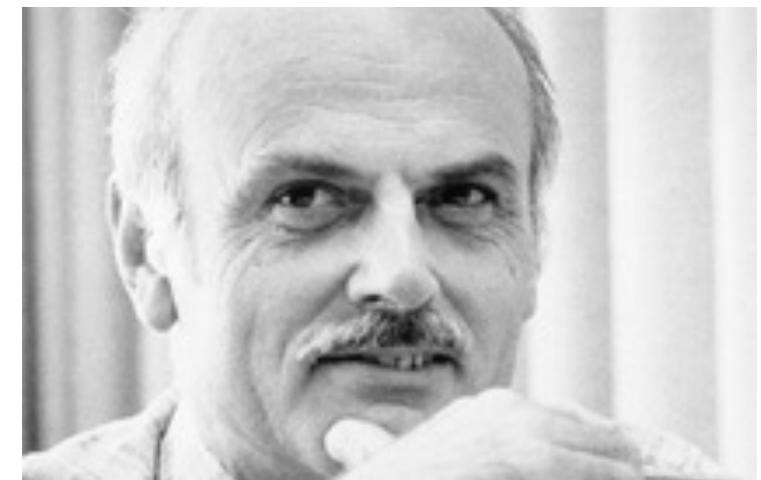
- Mid 1960s / early 1970s
- Hard disks allow for a **direct access** to data
- The time to locate a record is now measured in milliseconds
- New, more complex data structures were developed, such as lists and trees to be stored on disk
- Along with the direct access to data, it came a new type of system software known as a **database management system** (DBMS)



- With DBMS it was easier to access and manage data on a disk; moreover, the DBMS took care of tasks such as storing data, indexing data, governing access rights, and so forth
- By the mid-1970s, **online transaction processing** (OLTP) made even faster access to data possible
- **Interactive applications** like bank teller systems and manufacturing control systems became more and more widespread
- At this point, network and hierarchical data models were of widespread use

Meanwhile, Ted Codd defines the **relational data model**:

- Would win the ACM Turing Award for this work
- Points of strength: simplicity and hiding of physical details
- IBM Research begins work on *System R* prototype
- UC Berkeley begins work on *Ingres* prototype
- They both still exhibit computationally worse performance with respect to previous network and hierarchical data models



1980s:

- Relational DB prototypes evolve into commercial systems, overthrowing previous data models
- SQL becomes an industrial standard: programmers can now work at the logical level, without worrying about query optimality and the actual location of data
- Parallel and distributed database systems, on dedicated architectures
- Object-oriented database systems: they use the same model as object-oriented programming languages

SQL

- Around 1990 Wal-Mart retail corporation began to achieve wide acclaim for its mastery of supply chain management
- Behind this success there was a **data warehouse**, and a new way of interacting with data, called **online analytical processing** (OLAP)
- Data is collected by its point-of-sales systems to achieve unprecedented insight into the purchasing habits of its 100 million customers and the logistics guiding its 25,000 suppliers
- Wal-Mart's data warehouse was the first commercial Enterprise Data Warehouse to reach 1 terabyte of data in 1992



- Around 2000s, the types of data stored in database systems evolved rapidly, pushed by an ever increasing usage of the Internet and multimedia (Big Data)
- The variety of new data-intensive applications led to **NoSQL** systems, which gave programmers greater flexibility to work with new types of data, but lacked a high level query language
- **Distributed** storage and computing framework, capable of running on commodity hardware, were developed, such as Hadoop
- To allow for the interchange of information between systems, formats such as XML and JSON became of widespread usage
- Data Mining applications started to emerge
- Most of large and medium sized organizations started to rely on **decision support systems**

- A decision is the selection of a course of action from a set of alternatives
- A decision support system recommends such an action by offering managers information upon which to build ideas so to come up with the final choice
- A **decision management system** is an *action-oriented* evolution of a decision support system
- It makes one step more and takes decisions without human intervention based on known information and a set of coded business rules or **artificial intelligence models**
- Of course, not all judgments may be automated (strategic vs operational decisions)

- From 2010s, more and more companies began to shift from an on-premise management of their decision support / management systems to **cloud solutions** (e.g., offered by Microsoft or Amazon)
- Another important revolution in the field, started in 2023, is the advent of **Large Language Models** (LLMs), that promise support to the design and exploitation of data management systems



- A. Silberschatz, H.F. Korth, S. Sudarshan, *Database system concepts*, 7th Edition, 2020
- W.H. Inmon, *Building the Data Warehouse*, 4th Edition, 2005