# [IF-5-OT7:TD] Foundation of data engineering

## MCF Riccardo Tommasini
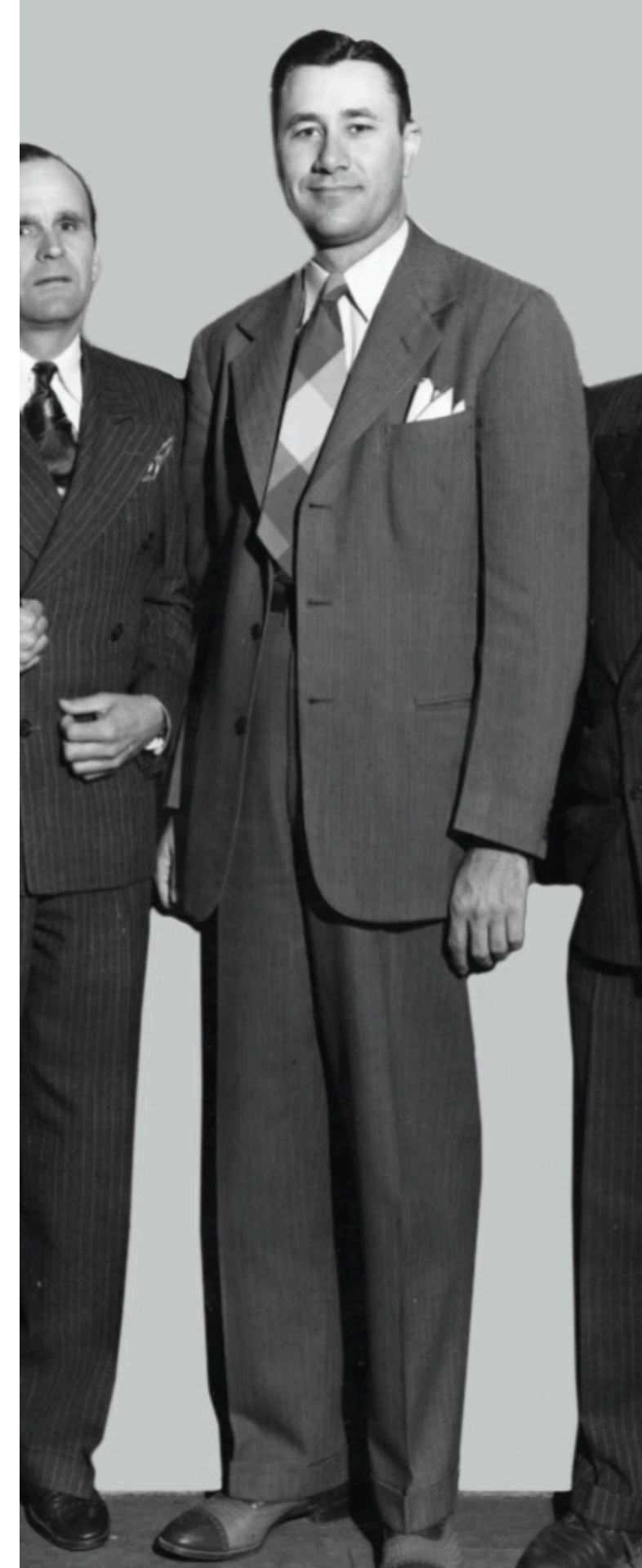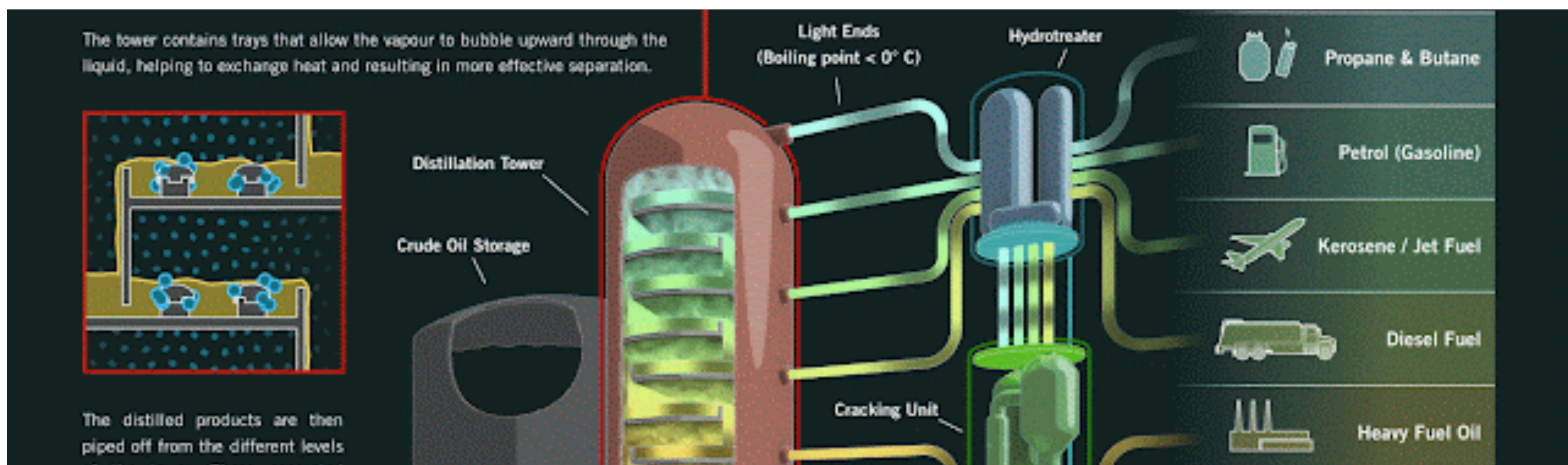
http://rictomm.me

riccardo.tommasini@insa-lyon.fr

# Quote

"A scientist can discover a new star, but he cannot make one.
He would have to ask an engineer to do it for him."

— *Gordon Lindsay Glegg*

# Data Science is...[01]



...refining crude oil

Discovering what we don't know from data
- Obtaining predictive, actionable insight from data
- Creating Data Products that have business impact now
- Communicating relevant business stories from data
- Building confidence in decisions that drive business value

# Data Engineering is...



...build the refinery.
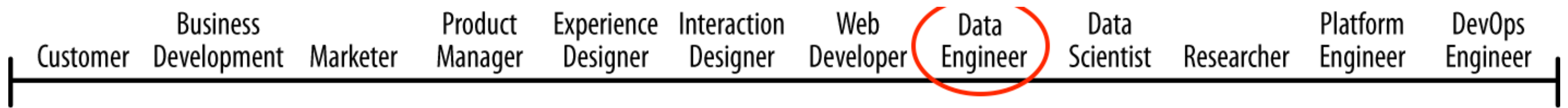
# Roles in a Data Science Project[02]

Customer | Business Development | Marketer | Product Manager | Experience Designer | Interaction Designer | Web Developer | Engineer | Data Scientist | Researcher | Platform Engineer | DevOps Engineer

---

# Roles in a Data Science Project[02]

Customer — Business Development — Marketer — Product Manager — Experience Designer — Interaction Designer — Web Developer — **Data Engineer** — Data Scientist — Researcher — Platform Engineer — DevOps Engineer
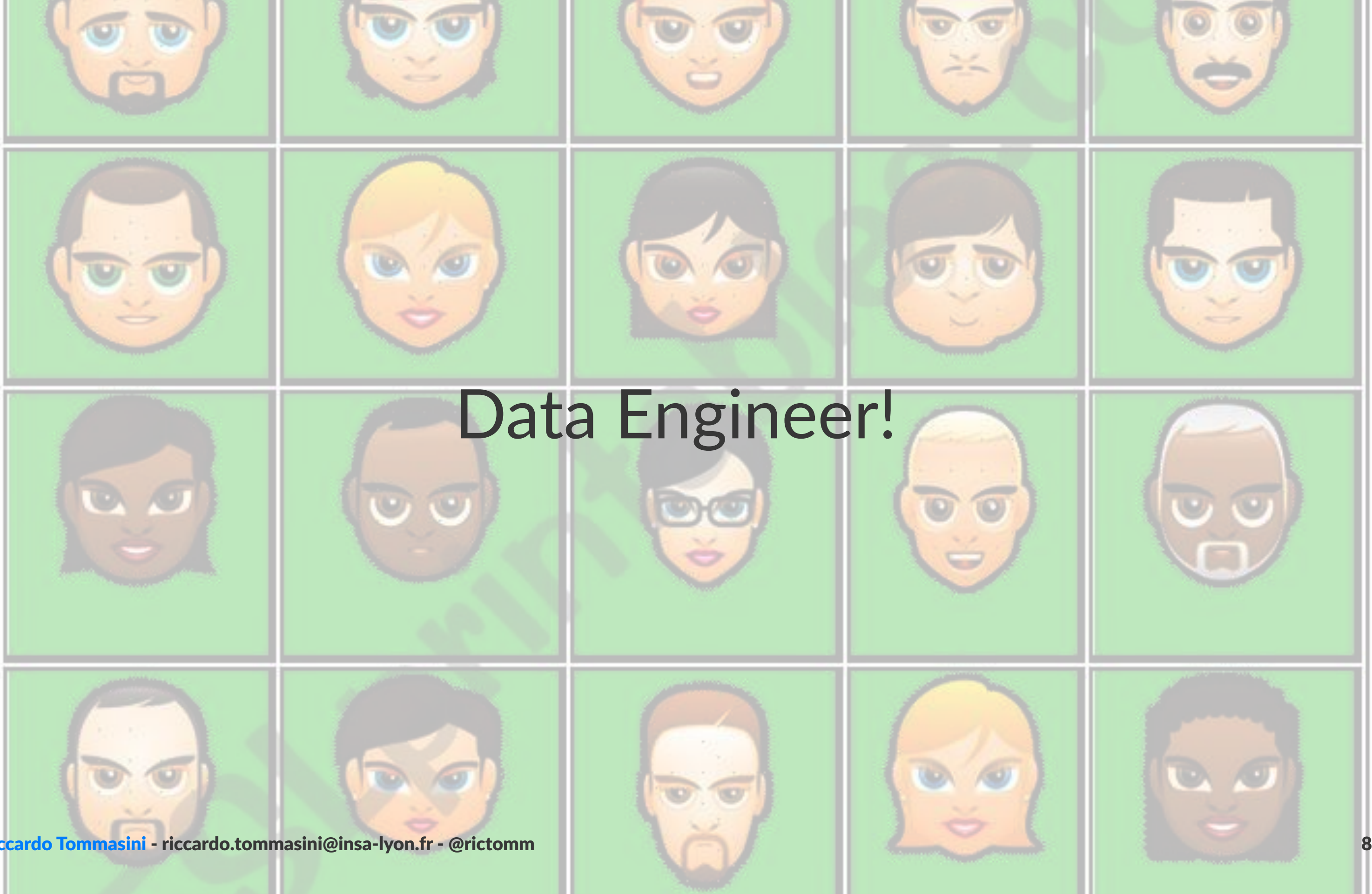
Data Engineer!

# The Data Engineer

A dedicated specialist that maintain data available and usable by others (Data Scientists).[03]

Data engineers set up and operate the organization's data infrastructure preparing it for further analysis by data analysts and scientists.[03]
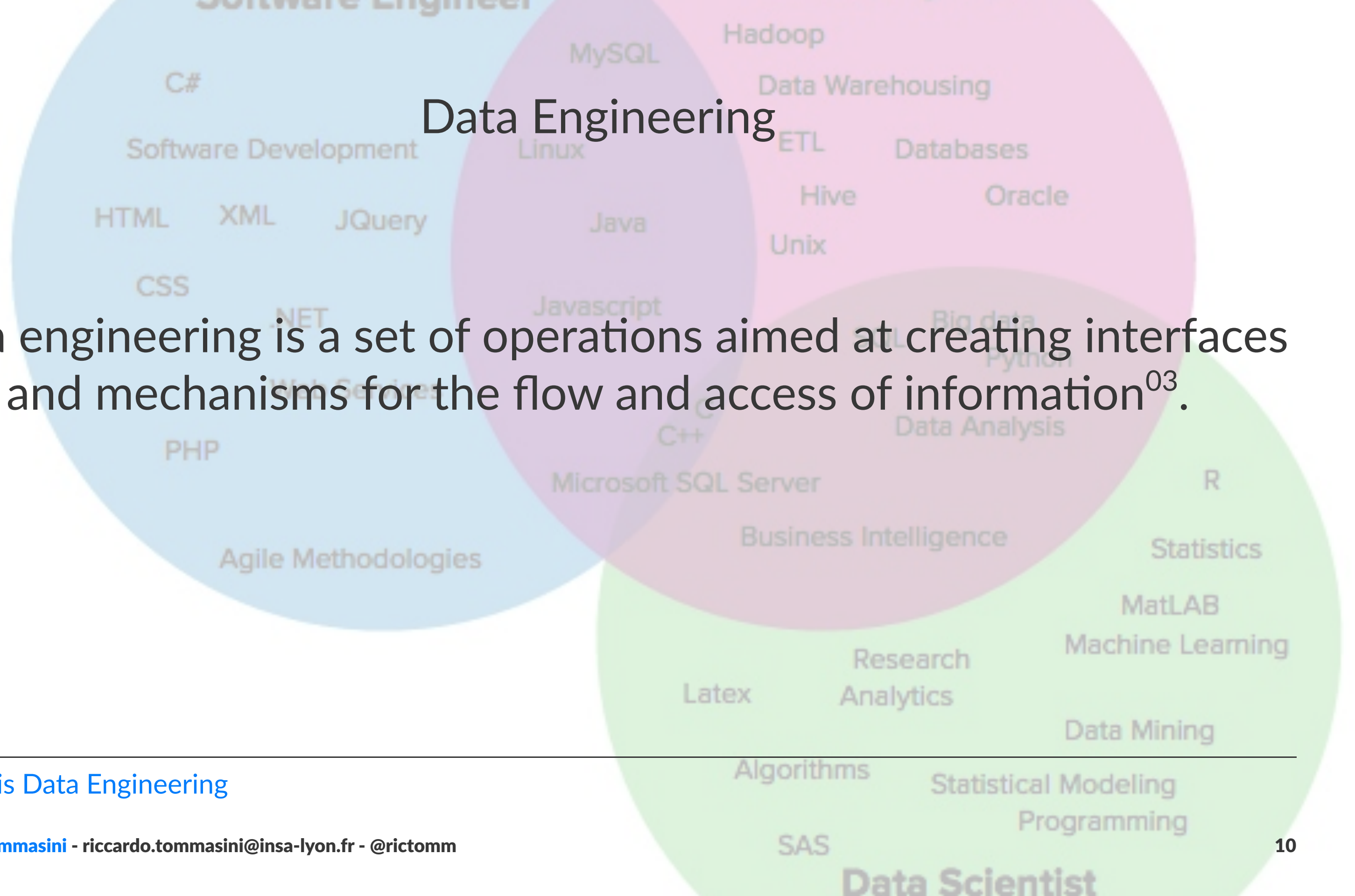
Data engineering field could be thought of as a superset of business intelligence and data warehousing that brings more elements from software engineering.[04]

---

[03] What is Data Engineering
[04] Source: The Rise of Data Engineer

# Data Engineering

Data engineering is a set of operations aimed at creating interfaces and mechanisms for the flow and access of information[03].

---

[03] What is Data Engineering

Venn diagram showing overlapping skill sets for Software Engineer, Data Engineer, and Data Scientist.

**Software Engineer**
- C#
- Software Development
- HTML
- XML
- JQuery
- CSS
- .NET
- Web Services
- PHP
- Agile Methodologies

*Software Engineer ∩ Data Engineer*
- MySQL
- Linux
- Java
- Javascript

**Data Engineer**
- Hadoop
- Data Warehousing
- ETL
- Databases
- Hive
- Oracle
- Unix

*Software Engineer ∩ Data Engineer ∩ Data Scientist*
- C
- C++
- Microsoft SQL Server

*Data Engineer ∩ Data Scientist*
- SQL
- Big data
- Python
- Data Analysis
- Business Intelligence

**Data Scientist**
- R
- Statistics
- MatLAB
- Machine Learning
- Research
- Analytics
- Latex
- Data Mining
- Algorithms
- Statistical Modeling
- Programming
- SAS

- a data engineer might create a new aggregate of a dataset containing trillions of streaming events

- analytics engineer might use that aggregate in a new report on global streaming quality

- a data scientist might build a new streaming compression model reading the report

---

[05] Netflix Innovation

# each of these workflows has multiple overlapping tasks:

# Data Bricks

# Google's Two-Cents

## Professional Data Engineer

A Professional Data Engineer enables data-driven decision making by collecting, transforming, and publishing data. A Data Engineer should be able to design, build, operationalize, secure, and monitor data processing systems with a particular emphasis on security and compliance; scalability and efficiency; reliability and fidelity; and flexibility and portability. A Data Engineer should also be able to leverage, deploy, and continuously train pre-existing machine learning models.

The Professional Data Engineer exam assesses your ability to:

✓ Design data processing systems

✓ Build and operationalize data processing systems

✓ Operationalize machine learning models

✓ Ensure solution quality

[Register]  [FAQs]

This exam is available in English and Japanese.

# The Knowledge Scientist<sup>06</sup>



---

---

# Nowdays we deal with a number of data from different domains.

# What is Data?

# Oxford Dictionary

*Data [**uncountable, plural**] facts or information, especially when examined and used to find out things or to make decisions.* [08]
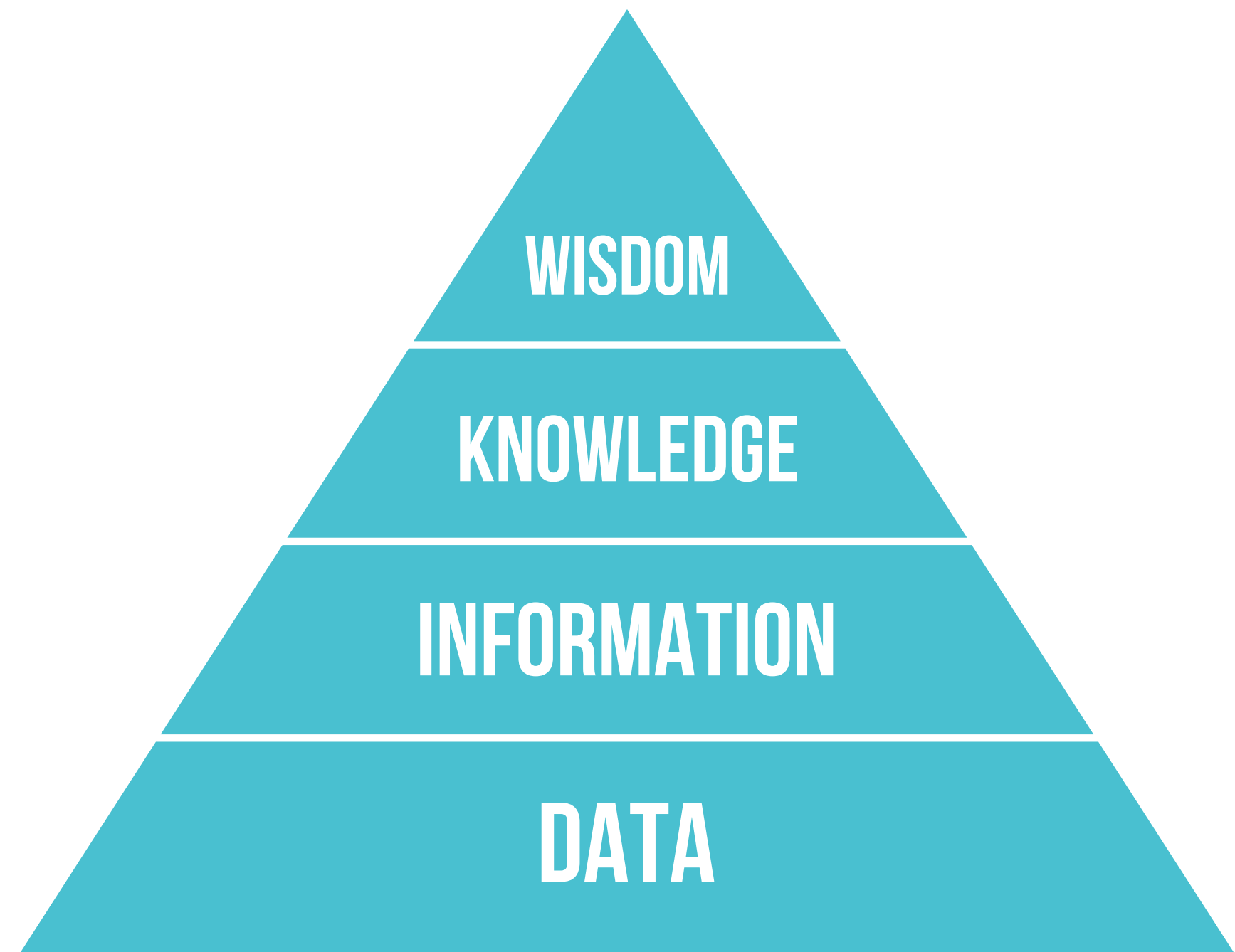
---

[08] Def

# Wikipedia

Data (treated as singular, plural, or as a mass noun) is any sequence of one or more symbols given meaning by specific act(s) of interpretation [09]

---

[09] Data in Computing)

# DIKW Pyramid

# Graph View

Data about data



METADATA!

# Data Semantics

## semantics
/sɪˈmantɪks/ 🔊

*noun*

the branch of linguistics and logic concerned with meaning. The two main areas are *logical semantics*, concerned with matters such as sense and reference and presupposition and implication, and *lexical semantics*, concerned with the analysis of word meanings and relations between them.

- the meaning of a word, phrase, or text.
  plural noun: **semantics**
  "such quibbling over semantics may seem petty stuff"

∨ Translations, word origin, and more definitions

# Big Data

# Challenges [014]



Data Velocity

Data Volume

Data Variety

real time
near real time
periodic
batch
table
data base
audio
unstructured
mobile
social
photo
web
video

MB   GB   TB   PB

Source: German Informatics Society

---

[014] Lanely, 2001

# Paradigm Shift

Big Data & Analytics

IBM

Paradigm shifts enabled by big data
Reduce effort required to leverage data

TRADITIONAL APPROACH

Small amount of carefully organized information

Carefully cleanse information
*before* any analysis

BIG DATA APPROACH

Large amount of messy information

Analyze information as is,
cleanse as needed

# Paradigm shifts enabled by big data
## Data leads the way—and sometimes correlations are good enough

**TRADITIONAL APPROACH**



Hypothesis → Question
Answer ← Data

Start with hypothesis and
test against selected data

**BIG DATA APPROACH**



Data → Exploration
Insight ← Correlation

Explore *all* data and
identify correlations

13

© 2014 IBM Corporation

# Paradigm shifts enabled by big data
## Leverage data as it is captured



**TRADITIONAL APPROACH**

Data → Repository → Analysis → Insight

Analyze data *after* it's been processed and landed in a warehouse or mart

**BIG DATA APPROACH**

Data
Analysis
Insight

Analyze data *in motion* as it's generated, in real-time

**Riccardo Tommasini** - riccardo.tommasini@insa-lyon.fr - @rictomm

# New Roles

In the context of Big Data, a data engineer must focus on **distributed systems**, and **programming languages** such as Java and Scala.
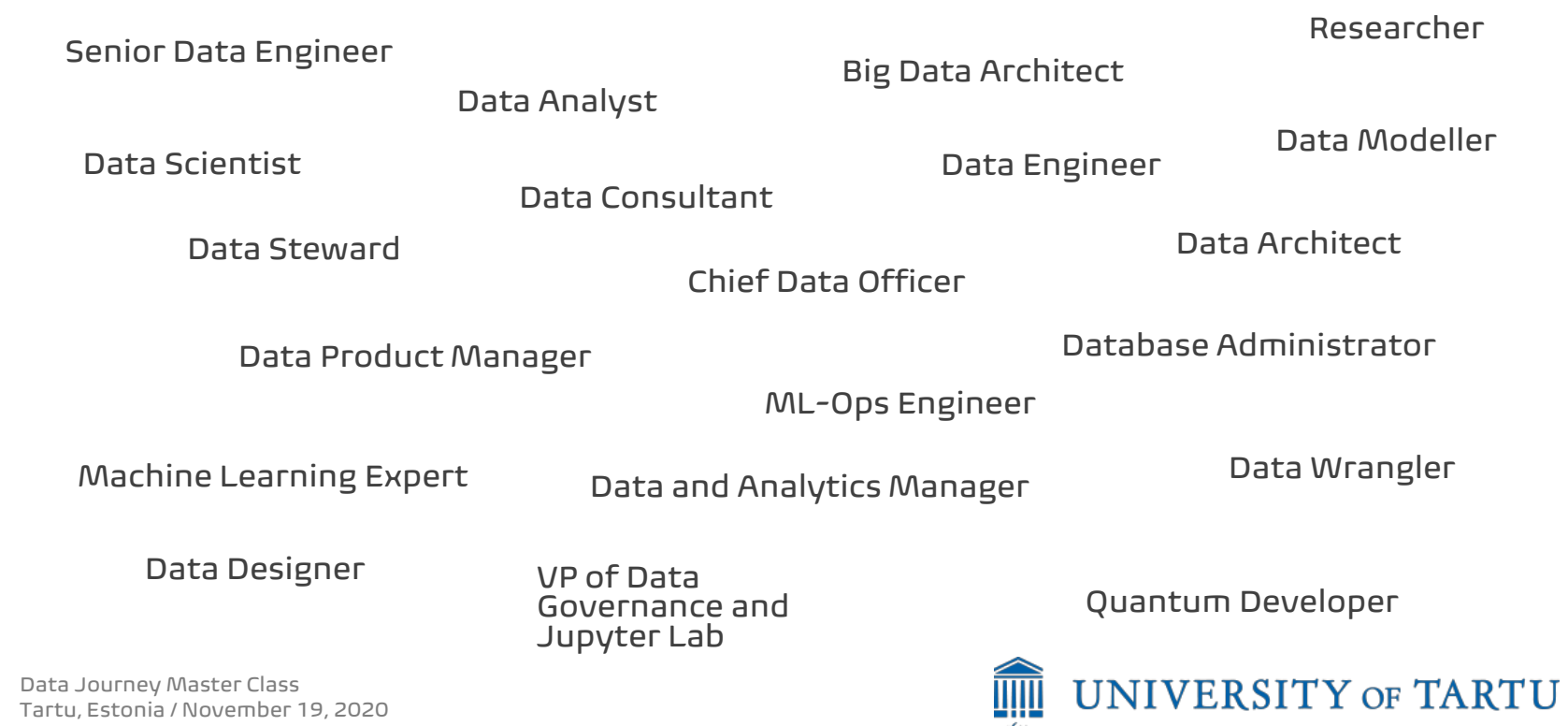
## Profession

Senior Data Engineer

Big Data Architect

Researcher

Data Analyst

Data Scientist

Data Engineer

Data Modeller

Data Consultant

Data Steward

Data Architect

Chief Data Officer

Data Product Manager

Database Administrator

ML-Ops Engineer

Machine Learning Expert

Data and Analytics Manager

Data Wrangler

Data Designer

VP of Data Governance and Jupyter Lab

Quantum Developer

Data Journey Master Class
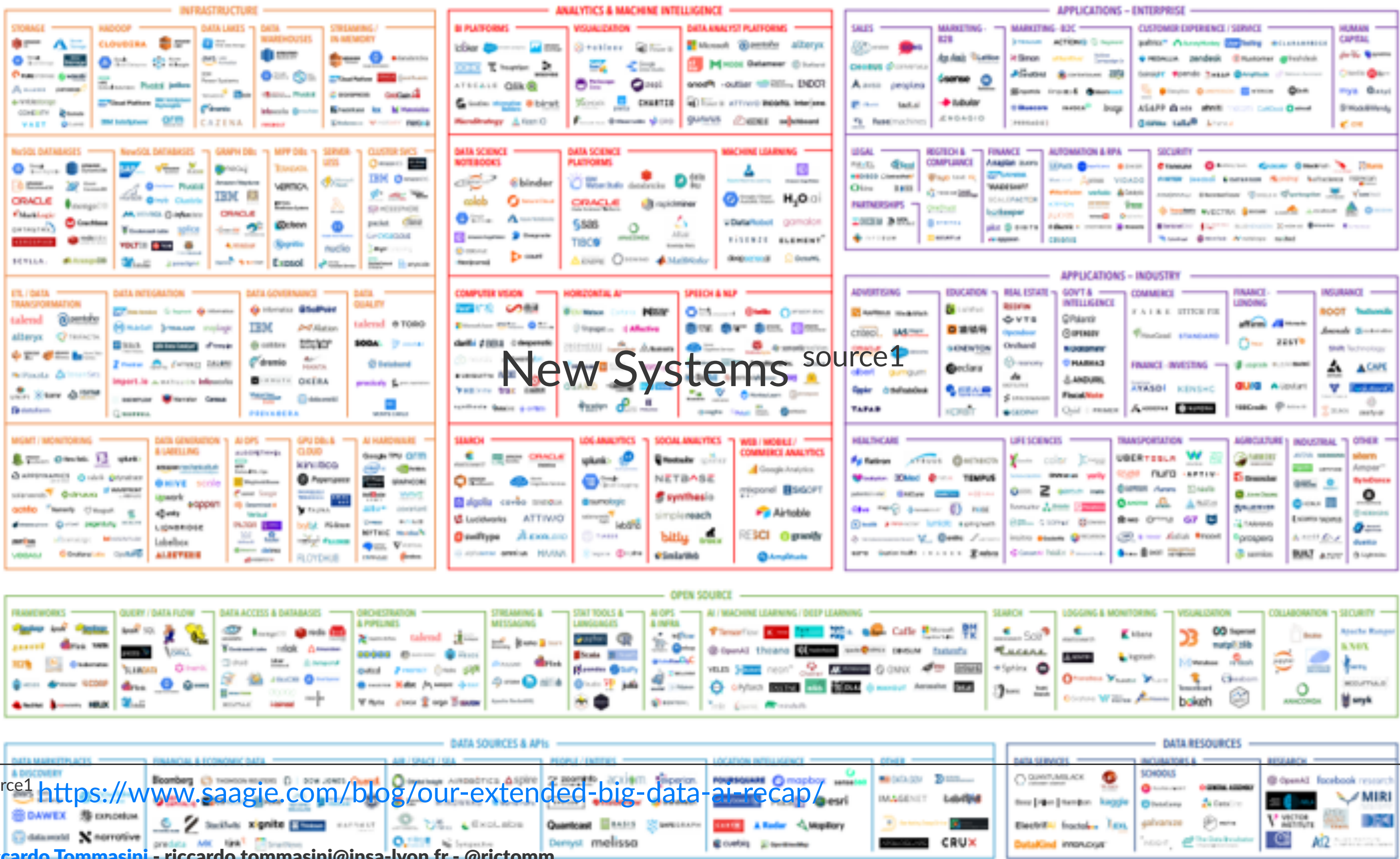Tartu, Estonia / November 19, 2020

UNIVERSITY OF TARTU

# New Tasks

Since data lake are taking data from a wide range of systems, data can be in **structured** or **unstructured** formats, and usually **not clean**, e.g., with missing fields, mismatched data types, and other data-related issues.

Therefore data engineers are challenged with the task of wrangling, cleansing, and integrating data.

New Systems    source1

# Where is Data?

# Data on the Inside vs Data on the Outside [pt]

| | Outside Data | Inside Data |
|---|---|---|
| Immutable? | Yes | No |
| Identity-Based References | Yes | No |
| Open Schema? | Yes | No |
| Represent in XML? | Yes | No |
| Encapsulation Useful? | No | Yes |
| Long-Lived Evolving Data with Evolving Schema? | No | Yes |
| Business Intelligence Desirable over Data? | Yes | Yes |
| Durable Storage in SQL Inside the Service? | Yes: Copy of XML Kept in SQL | Yes |

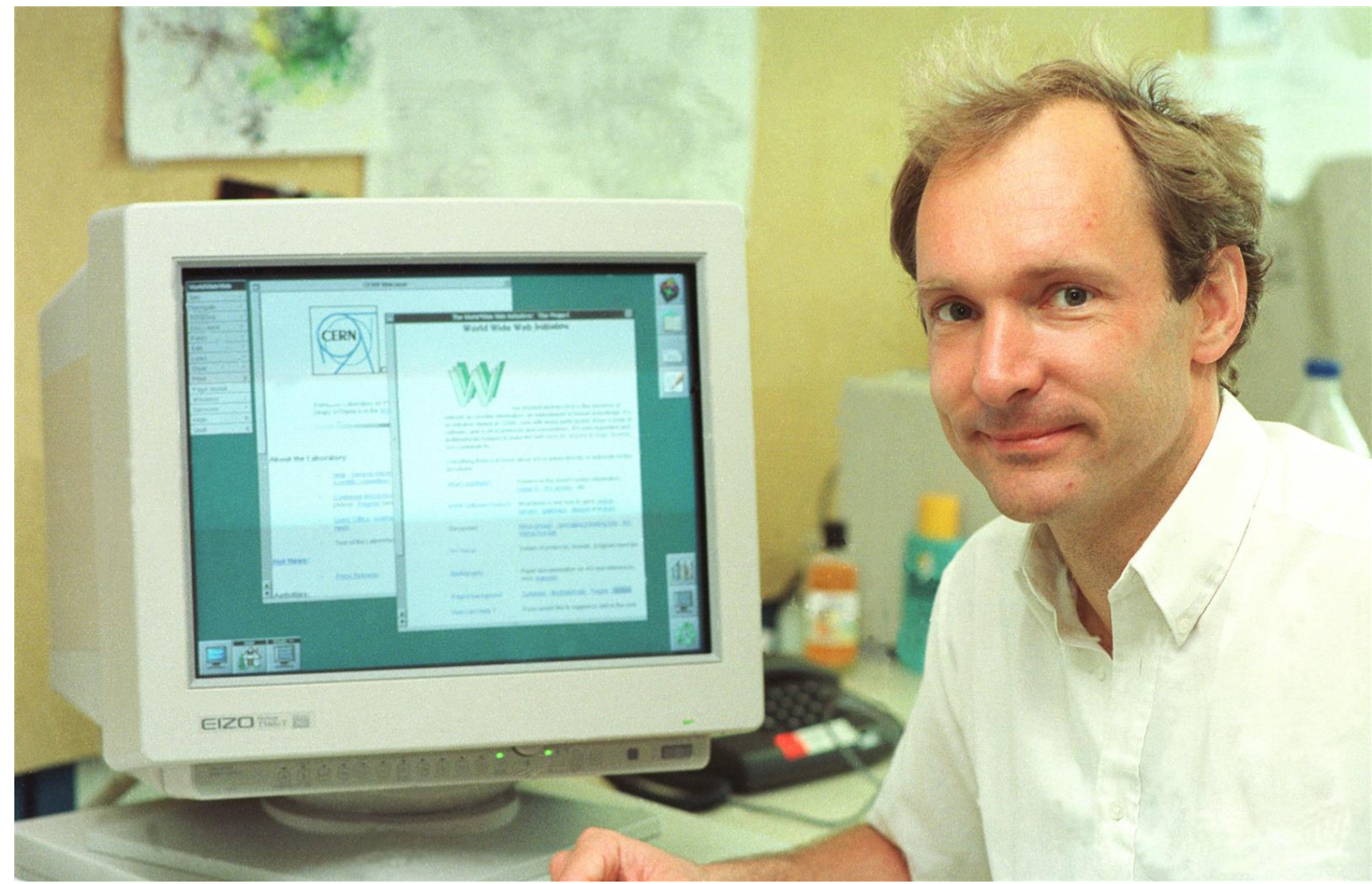[pt] Data on the Outside vs Data on the Inside *Pat Helland, CIDR 2005*

**Riccardo Tommasini** - riccardo.tommasini@insa-lyon.fr - @rictomm

37

# The Outside (WEB)

International ecosystem of applications and services that allows us to search, aggregate, combine, transform, replicate, cache, and archive the information that underpins society.

The Web is the result of millions of simple, small-scale interactions between agents and resources that use the founding technologies of HTTP and URI. [121]

The Web is a set of widely commoditised servers, proxies, caches, and content delivery networks [an engineers]



---

[121] Architecture of the World Wide Web, Volume One

**Riccardo Tommasini** - riccardo.tommasini@insa-lyon.fr - @rictomm

# Resources

Resources are the fundamental building blocks

Anything we can expose, i.e., documents, images, videos, audio, devices, people, things…

We can represent them by abstracting the useful information and identifying using a Uniform Resource Identifier (URI)

# URIs

URL format (RFC 2396):
scheme:[//[user:password@]host[:port]][/]path[?query][#fragment]



E.G.:
git@github.com:nodejs/node.git
mongodb://root:pass@localhost:27017/TestDB?options?replicaSet=test
http://example.com

# Representation

Access to a resource is mediated by a representation

This separation is convenient to promote loose-coupling between server (producers) and client (consumers)

Multiple Views and Content Negotiation are the basis for interoperability

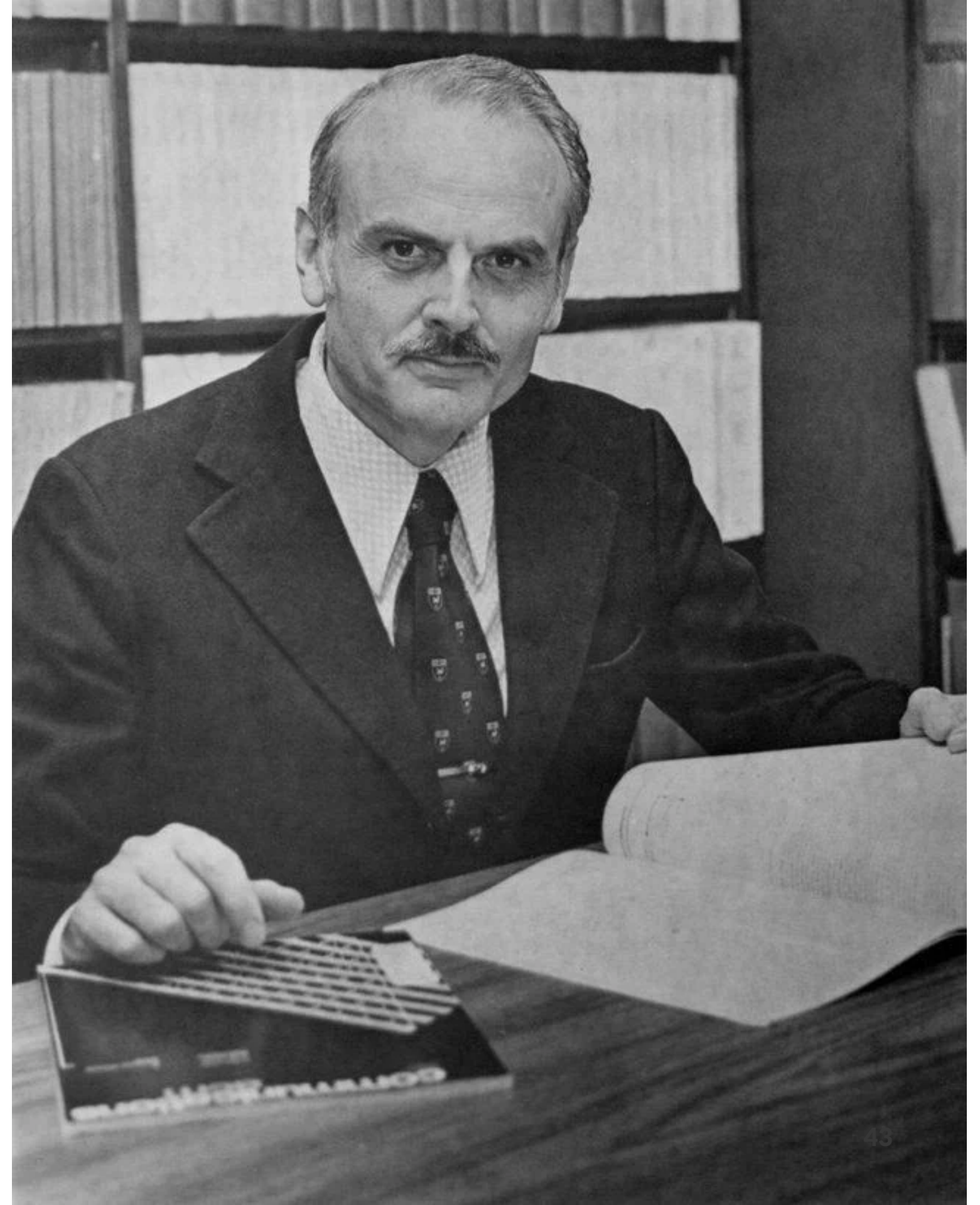# Protocols: HTTP

- GET

  - Uniform Interface

  - read-only operation

  - idempotent

- POST

  - like a resource upload

  - idempotent

- DELETE

  - remove resources

  - idempotent

- HEAD

  - HEAD is like GET except it returns only a response code

- PUT
  the only non-idempotent and unsafe operation
  is allowed to modify the service in a unique way

- OPTIONS is used to request information about the communication options of the resource

# Databases Management Systems

A database is **an organised collection of structured information, or data**, typically stored electronically in a computer system

Several kind of DBSMs exist. We will survey some of them. It is interesting to know that Edgar F. Codd defined 12+1 rules that make a DBMS relational link

**Riccardo Tommasini** - riccardo.tommasini@insa-lyon.fr - @rictomm

# Relational DBMS

- It must be relational as a database and as a management system

- All data should be in table form

- All data should be accessible without ambiguity



Codd's 13 Rules

- Foundation Rule
- Information Rule
- Non-Subversion Rule
- Distribution Independence
- Guaranteed Access Rule
- Integrity Independence
- Systematic Treatment of NULL Values
- Logical Data Independence
- Active Online Catalog
- Physical Data Independence
- High-level Insertion, Updation and Deletion
- View Updating Rule
- Comprehensive Data sub-language Rule

# NoSQL Familty

| Document Database | Graph Databases |
|---|---|
| Couchbase<br>MarkLogic™<br>mongoDB | Neo4j<br>InfiniteGraph<br>The Distributed Graph Database |
| **Wide Column Stores** | **Key-Value Databases** |
| Cassandra<br>amazon DynamoDB<br>riak<br>AEROSPIKE | accumulo<br>HYPERTABLE inc<br>redis  APACHE HBASE<br>Amazon SimpleDB |

# Data Warehouse: A Traditional Approach:

A data warehouse is a copy of transaction data specifically structured for query and analysis. — Ralph Kimball

A data warehouse is a subject-oriented, integrated, time-variant and non-volatile collection of data in support of management's decision making process.-- Bill Inmon

- A data warehouse is a central repository where raw data is transformed and stored in query-able forms.[03]

- Data Warehouse are still relevant today and their maintenance is part of Data Engineers' resposibilities.

- The warehouse is created with structure and model first before the data is loaded and it is called schema-on-write.

---

[03] What is Data Engineering

# Data Warehouse vs Data Bases

Surprisingly, Data Warehouse isn't a regular database.

- A database normalizes data separating them into tables and avoiding redundancies

- It supports arbitrary workload and complex queries

- do not store multiple versions of data

- a Data Warehouse uses few tables to improve performance and analytic.

- a Data Warehouse allows simple queries

- supports versioning for complex analysis

# Data Lake

A Data lake is a vast pool of raw data (i.e., data as they are natively, unprocessed). A data lake stands out for its high agility as it isn't limited to a warehouse's fixed configuration[03].

---

[03] What is Data Engineering

**HOW DO DATA LAKES WORK?**

The concept can be compared to a water body, a lake, where water flows in, filling up a reservoir and flows out.

**1** The incoming flow represents multiple raw data archives ranging from emails, spreadsheets, social media content, etc.

**STRUCTURED DATA**
1. Information in rows and columns
2. Easily ordered and processed with data mining tools

**UNSTRUCTURED DATA**
1. Raw, unorganized data
2. Emails
3. PDF files
4. Images, video and audio
5. Social media tools

**2** The reservoir of water is a dataset, where you run analytics on all the data.

**3** The outflow of water is the analyzed data.

**4** Through this process, you are able to "sift" through all the data quickly to gain key business insights.

Full Inforgraphic

- In Data Lake, the raw data is loaded as-is, when the data is used it is given structure, and it is called schema-on-read.

- Data Lake gives engineers the ability to easily change.

- In practice, Data Lake is a commercial term so don't sweat it.

# Data Lake vs Data Warehouse

**DATA WAREHOUSE**

**DATA LAKE**

VS

- **Structured Data**
- **Schema On Write**
- **Data Pipelines: Extract-Transform-Load**
- **Processing Model: Batch**

- **Unstructured Data**
- **Schema on Read**
- **Data Pipelines: Extract-Load-Transform**
- **Processing Model: Streaming**