



Laboratoire d'Informatique et d'Automatique pour les Systèmes



# Designing eco-friendly query processors

Ladjel BELLATRECHE

LIAS/ISAE-ENSMA, Poitiers, France

[bellatreche@ensma.fr](mailto:bellatreche@ensma.fr)

<http://www.lias-lab.fr/members/bellatreche>



École Saisonnière en Intelligence Artificielle  
14 avril 2025, Strasbourg



# Big Thank to

Organisateurs  
(efficacité)



**AfIA**

Association française  
pour l'Intelligence Artificielle

Ecole Saisonnière en IA

Dr. Ahmed SAMET



Energy savings



**Query Optimisation**

Improvement of database  
performance



**Energy Efficiency**

Quantification & reduction of  
consumed energy  
of queries/jobs



**A decarbonized world**

Reduction of greenhouse  
gas emissions

# ISAE-ENSMA

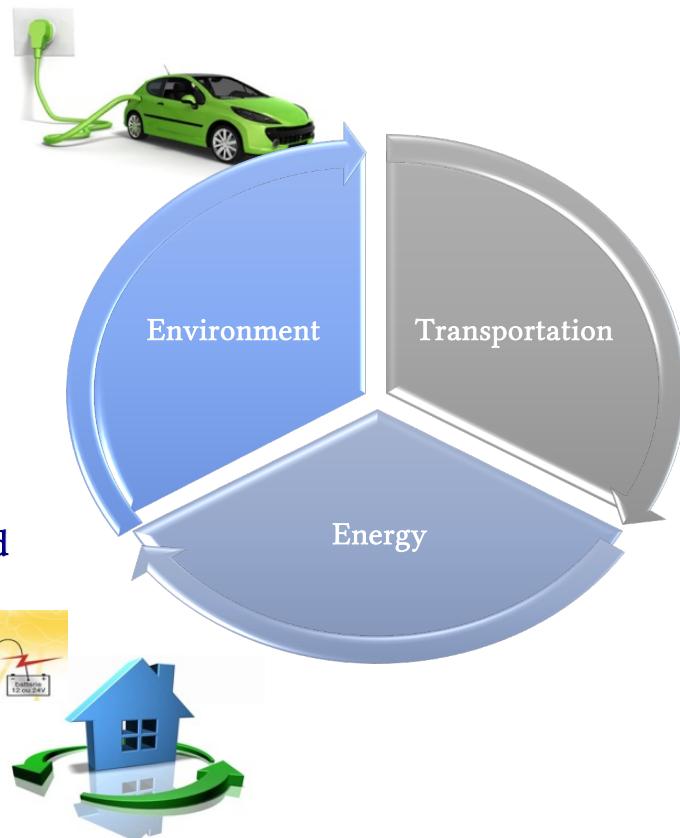
- ▶ National Engineering School for **Mechanics** and **Aerotechnics**, Futuroscope, Poitiers
- ▶ Member of ISAE-Group: SUPAERO, ENSMA, SUPMECA, ESTACA,, ENAC, Ecole de l'Air et de l'Espace
- ▶ Ranked by Shanghai Ranking (between 151-200)



# LIAS Laboratory

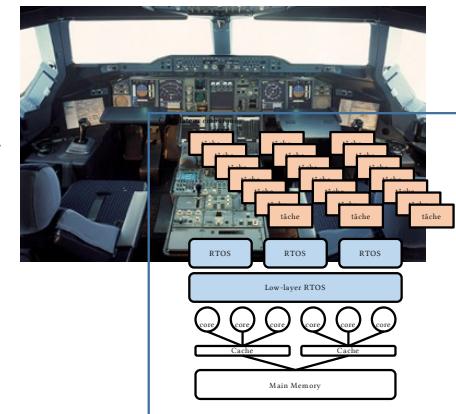
## Automation & Systems

- Modelling
- Analysis
- Command



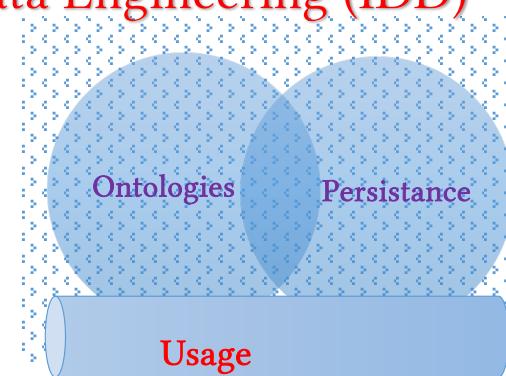
## Embedded & Real Time Systems

- Design
- Validation
- Sizing

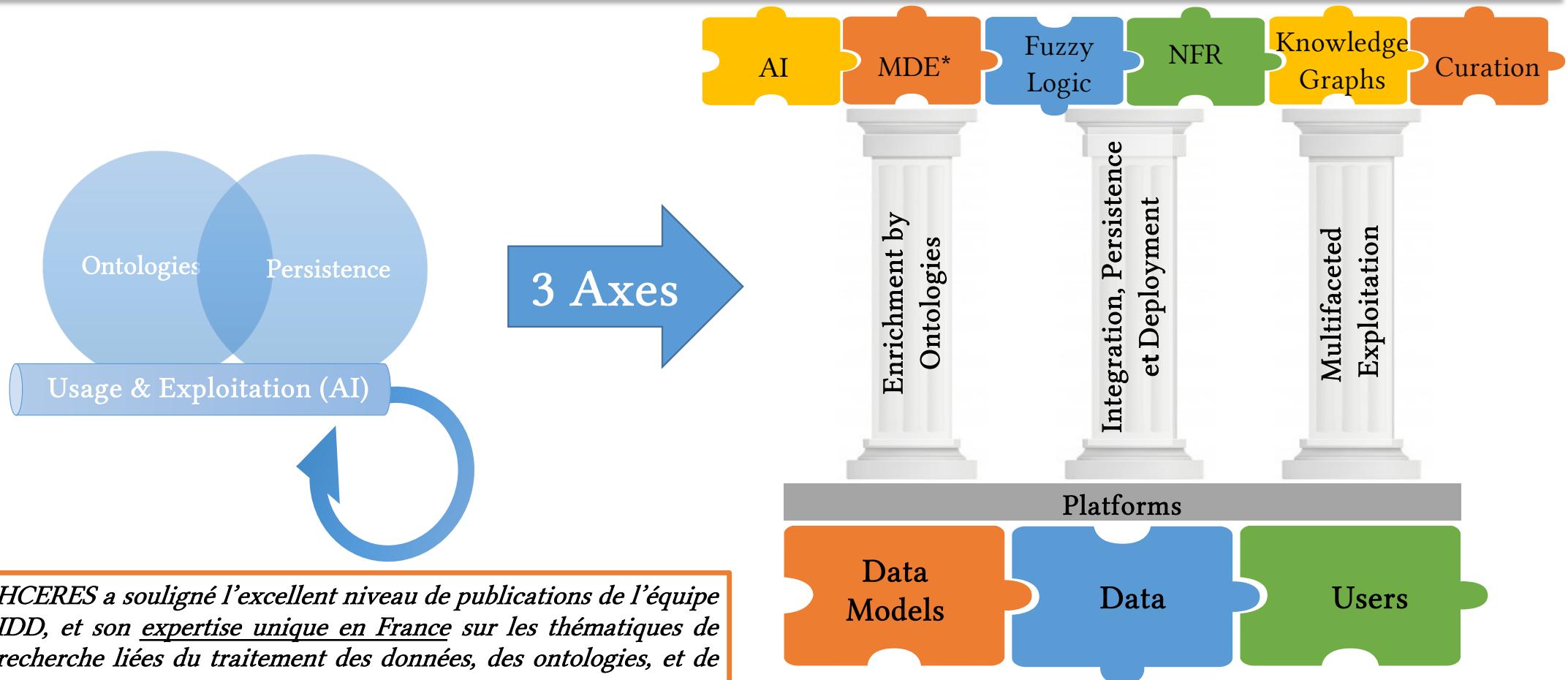


## Model & Data Engineering (IDD)

- Ontologies
- Persistence
- Usages



# Themes of IDD: [HCERES 2021]



MDE: Model Driven Engineering

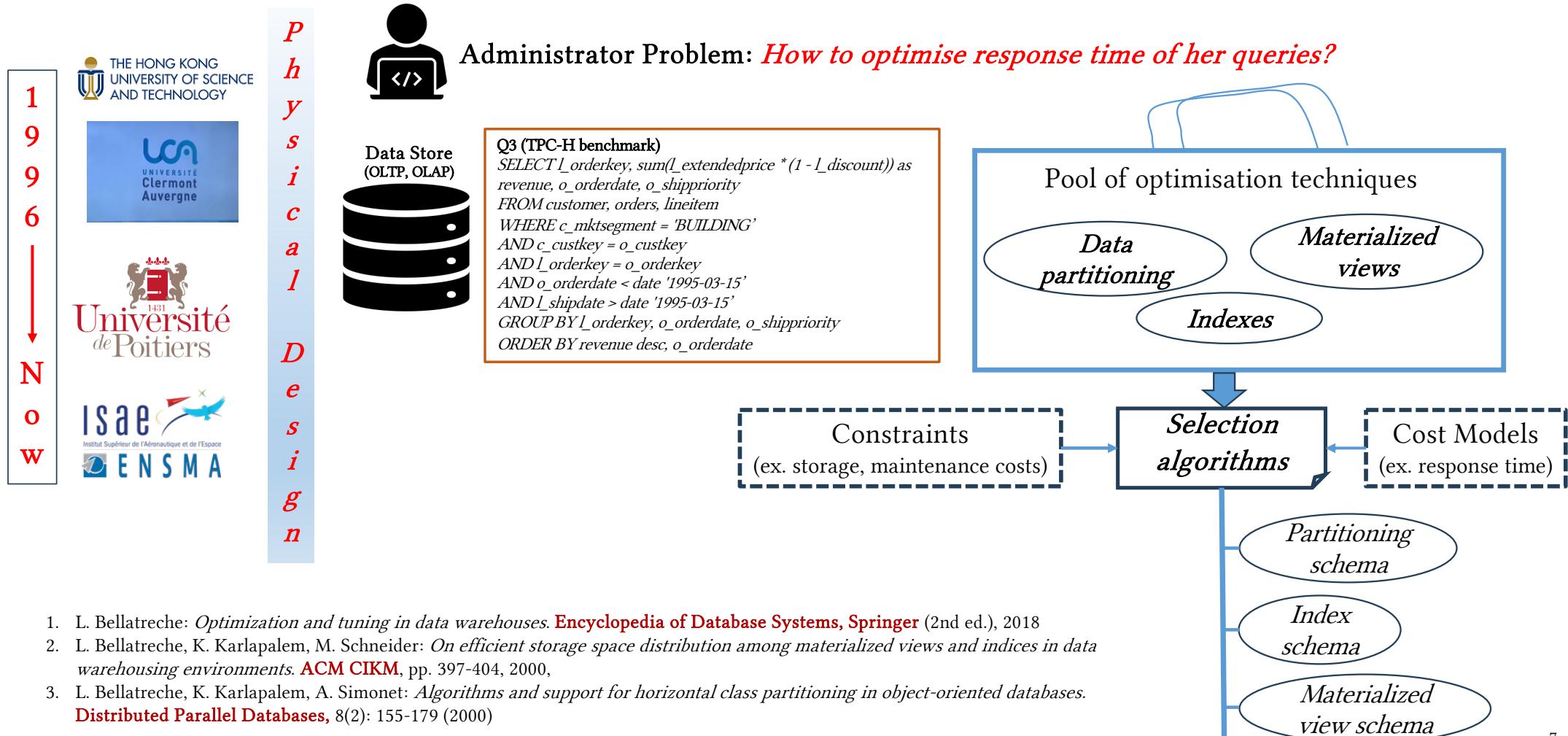
NFR: Non Functional Requirements (ex. performance, energy consumption)

# Agenda

---

1. My journey to “Green” Query Processors (QP)
2. Digitalisation, Data Science, and Energy
3. Framework for Studying Energy Efficiency
  - Energy Efficiency of QP
4. Summary

# My journey to “green” query processors



1. L. Bellatreche: *Optimization and tuning in data warehouses*. **Encyclopedia of Database Systems, Springer** (2nd ed.), 2018
  2. L. Bellatreche, K. Karlapalem, M. Schneider: *On efficient storage space distribution among materialized views and indices in data warehousing environments*. **ACM CIKM**, pp. 397-404, 2000,
  3. L. Bellatreche, K. Karlapalem, A. Simonet: *Algorithms and support for horizontal class partitioning in object-oriented databases*. **Distributed Parallel Databases**, 8(2): 155-179 (2000)

# Challenges of physical design

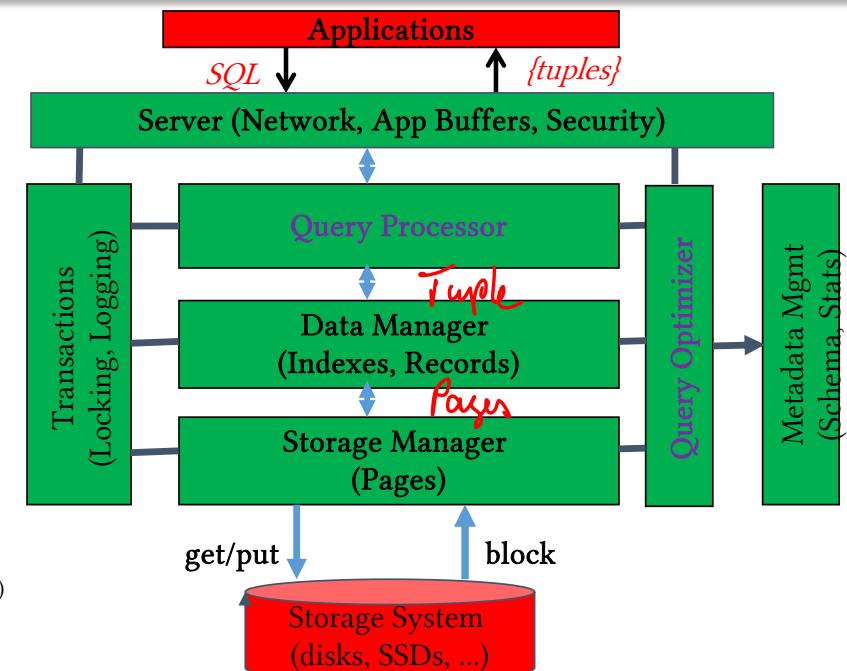
- ▶ Knowledge of functioning of a DBMS
- ▶ Definition an analytical cost model to estimate
  1. Response time of DB operations (e.g., join, selection)
  2. Maintenance cost
  3. Storage cost

- ▶ It concerns

- Elementary operation (ElemOp) and its implementation (ex. join, sorting)
- Metrics (MET): Inputs-Outputs (IO), CPU, and Network
- Involved parameters (P) (DBMS, database, query, optimization techniques, platform, hardware, ...)

$$CM_{ElemOp}^{MET_m}: P^n \rightarrow Value\ of\ MET_m \in \mathbb{R}$$

- ▶ Cost-based algorithms to select optimisation techniques (ex. genetic algorithms, simulated annealing, and hill climbing)



1. A. Ouared, Y. Ouhammou, L. Bellatreche: *QoS MOS: QoS metrics management tool suite*. **Computer Languages, Systems & Structures, Elsevier**. 54: 236-251 (2018)
2. A. Ouared, Y. Ouhammou, L. Bellatreche: *CostDL: A cost models description language for performance metrics in databases*. **ICECCS**, 187-190, 2016
3. S. Benkrid, Y. Mestoui, L. Bellatreche, C. Ordóñez: *A genetic optimization physical planner for Big Data warehouses*. **IEEE Big Data**, 406-412, 2020

# Cost model: an example

*Estimation of number of pages (IO) of: SELECT \* FROM T*

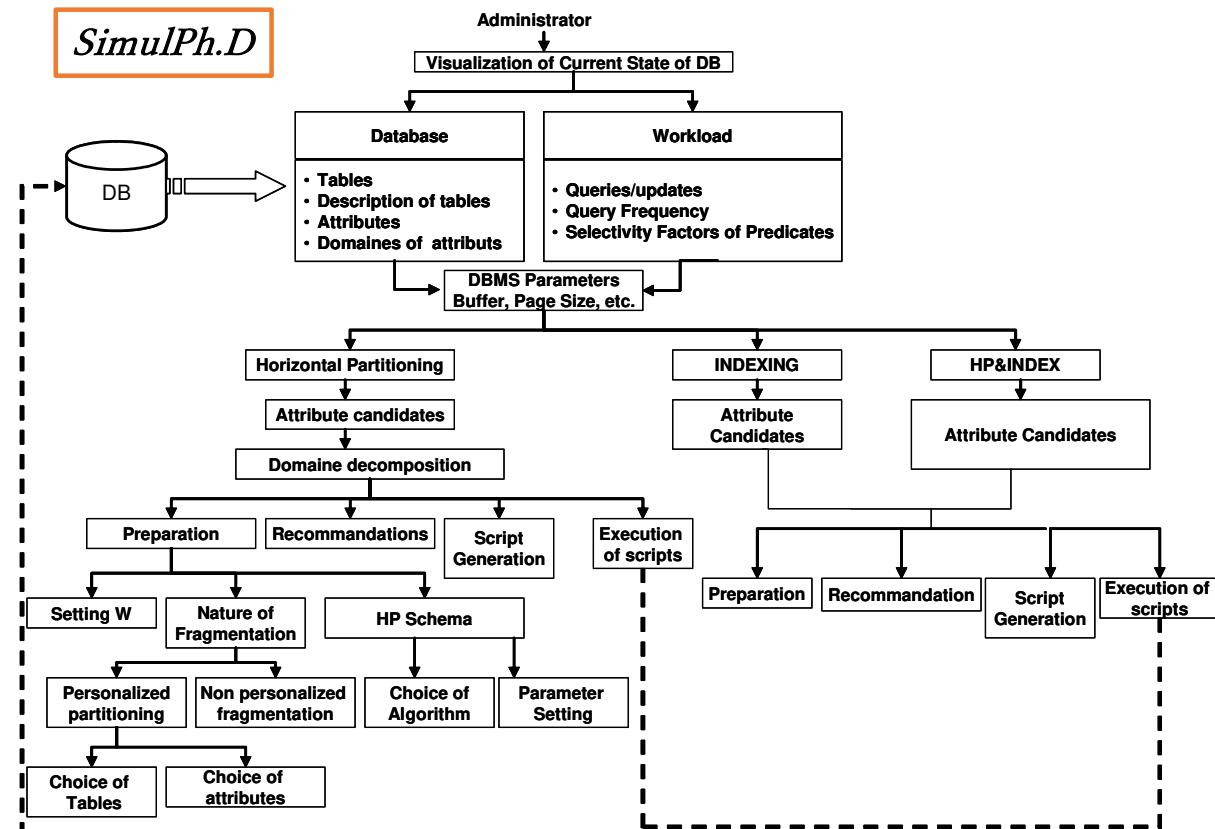
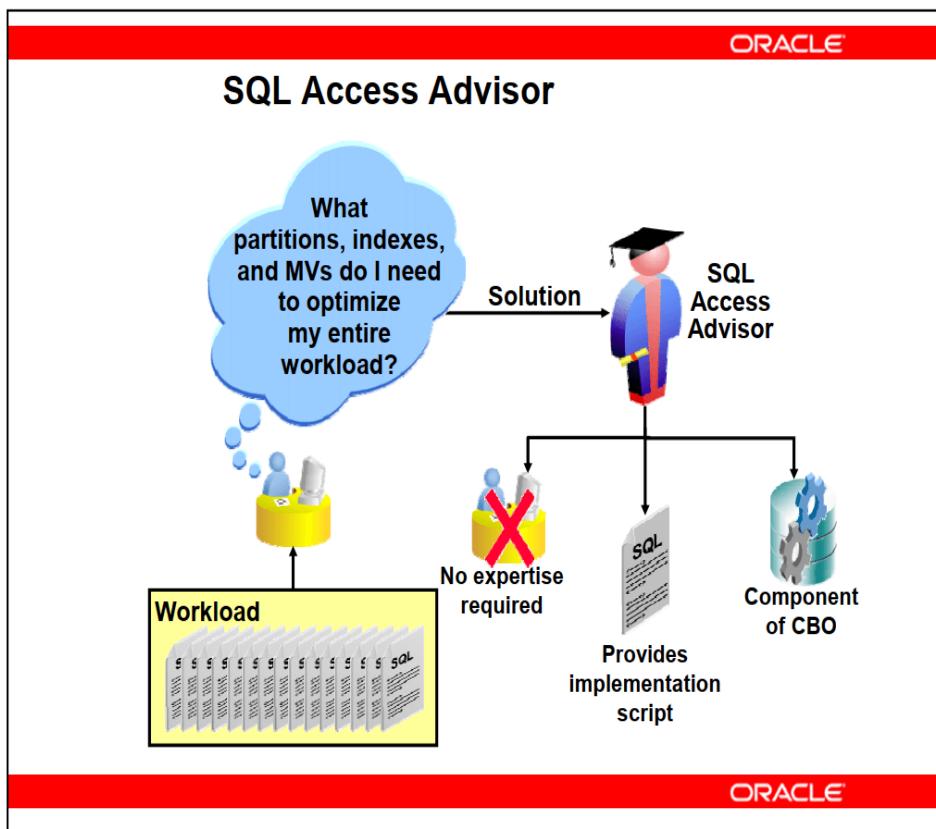
- ▶ Query optimizer functioning [row-store]: Disk-oriented DBMS
  - The CPU can only work on tuples that **are in memory**
  - Tuples must first be transferred from **disk to memory**
  - Data is transferred from disk to memory (and back) in whole blocks at the time
  - The disk can hold D blocks, at most M blocks can be in **buffer** at the same time ( $M \ll D$ )



$$CM_{\text{Projection}}^{\text{IO}} = \left\lceil \frac{\|T\| * LG}{\text{Bloc Size}} \right\rceil$$

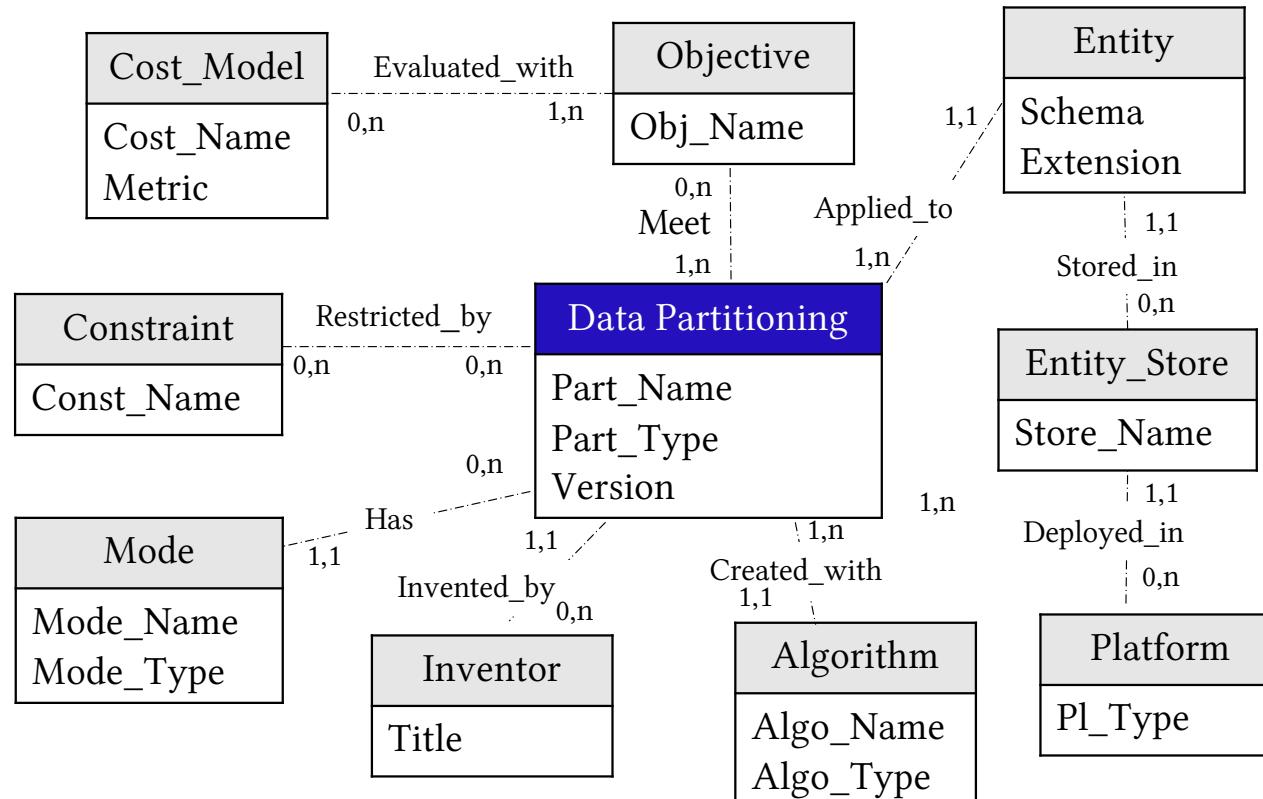
- $\|T\|$ : number of tuples of T
- LG: Length of a tuple of T

# Commercial and academic advisors



# Generalisation of our findings

→ The case of data partitioning for RDF data store



1. J. Galicia: *Revisiting data partitioning for scalable rdf graph processing*. **Ph.D. thesis**, ISAE-ENSMA, Poitiers, France, 2021
2. A. Khelil, A. Mesmoudi, J. Galicia, L. Bellatreche, M.-S. Hacid, E. Coquery: *Combining graph exploration and fragmentation for scalable RDF query processing*. **Information Systems Frontiers**, Elsevier, 23(1): 165-183, 2021

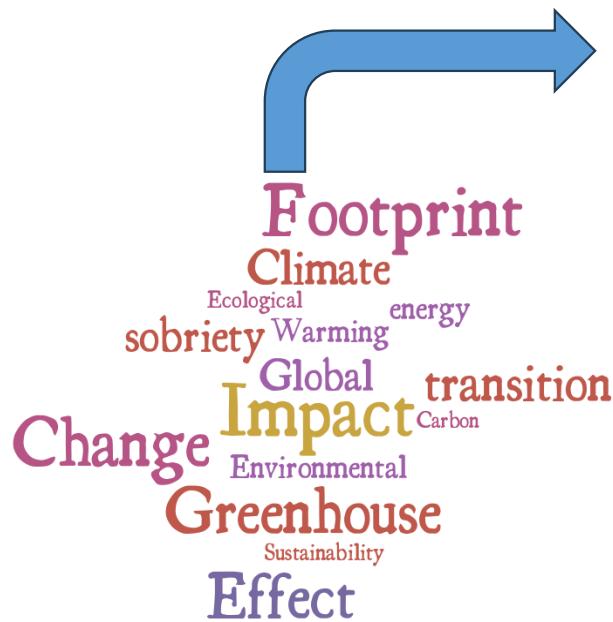
# Takeaways

- ▶ 11 Ph.D. theses (<https://theses.fr/112281834>)
- ▶ 20 masters



- A robust and well-established **framework** for physical design
- **Experience** in defining and refining **cost models**
- A diverse range of selection algorithms
- Collaboration with major DBMS companies (Teradata)

# 2015: COP21 @ Paris



Raising awareness on climate change (CC) and health

CC is everyone's business: green initiatives

Individuals

Scientists

National Institutions  
(Governments, Associations)

International Institutions  
(e.g., WHO)

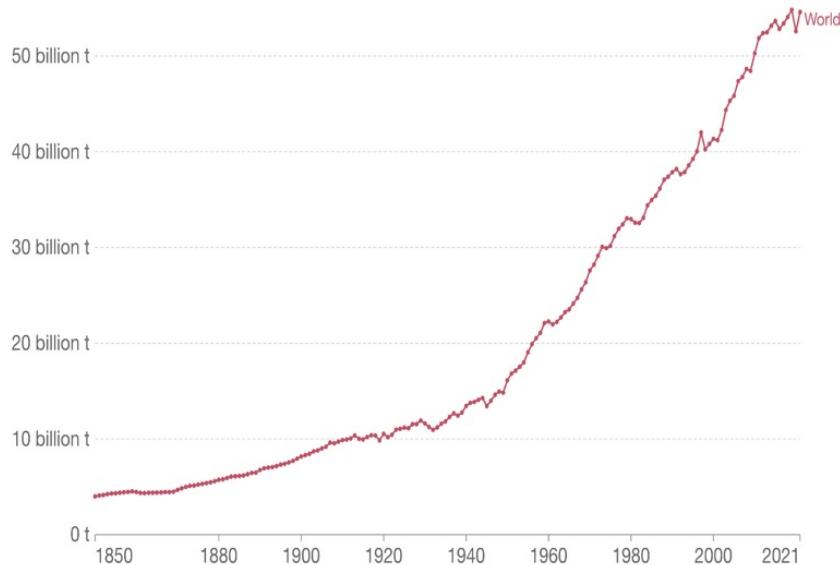
Green shoots in last French local elections  
(Lyon, Bordeaux, Poitiers, ...)

→ Energy efficiency (EE): achievement of the same level of services while consuming less energy



- ▶ How can I study EE of databases?
  - € Fundings (local, national, and EU)
  - 💻 A niche research topic

# Global greenhouse gas emissions



Source: Calculated by Our World in Data based on emissions data from Jones et al. (2023)

Note: Land use change emissions can be negative.

OurWorldInData.org/co2-and-greenhouse-gas-emissions • CC BY



## ► Main causes of global warming:

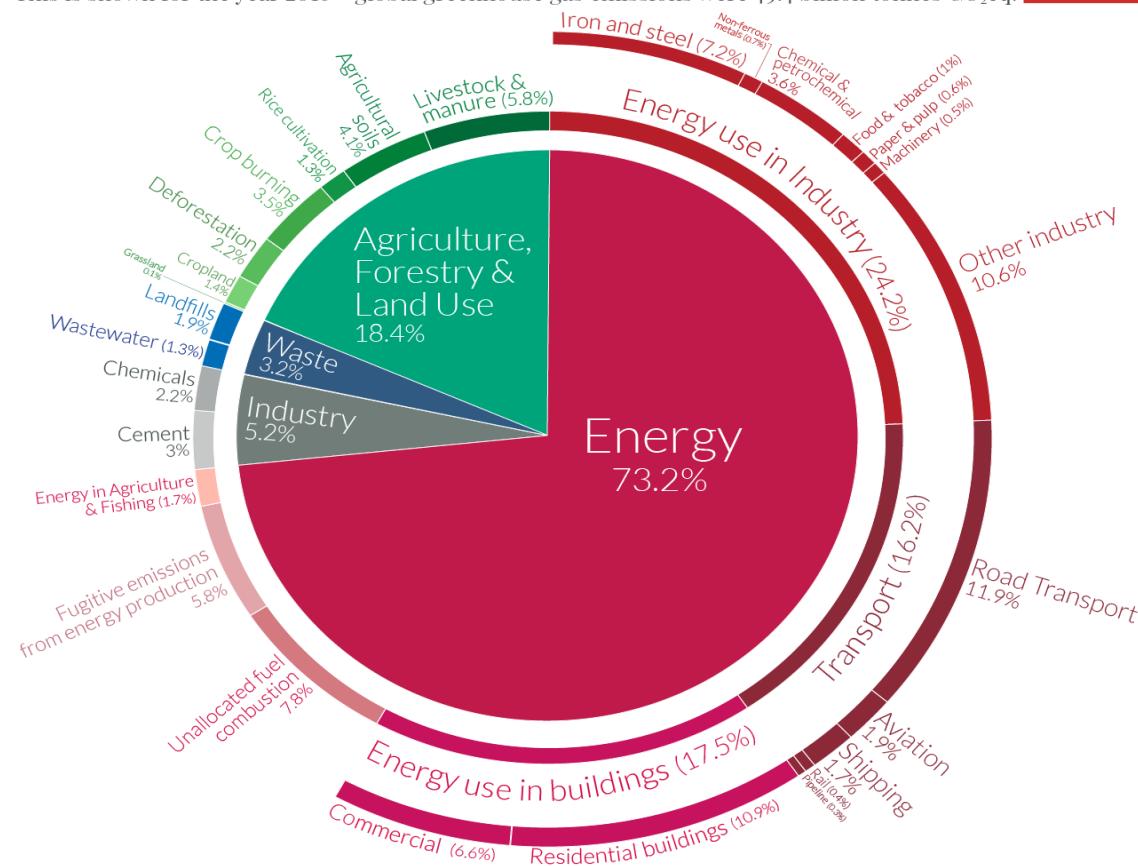
- Power Generation by burning fossil fuels
- Manufacturing industry is one of the largest contributors to greenhouse gas emissions worldwide
- Cutting down forests reduces the absorption of carbon dioxide
- Transportations that run on fossil fuels
- Production of Food
- **Powering buildings that consume over half of all electricity: electricity consumption for lighting, appliances, and connected devices**
- Our consumption lifestyle and high energy demand

<https://www.un.org/en/climatechange/science/causes-effects-climate-change>

# Greenhouse gas emission by sector

Global greenhouse gas emissions by sector  
This is shown for the year 2016 – global greenhouse gas emissions were 49.4 billion tonnes CO<sub>2</sub>eq.

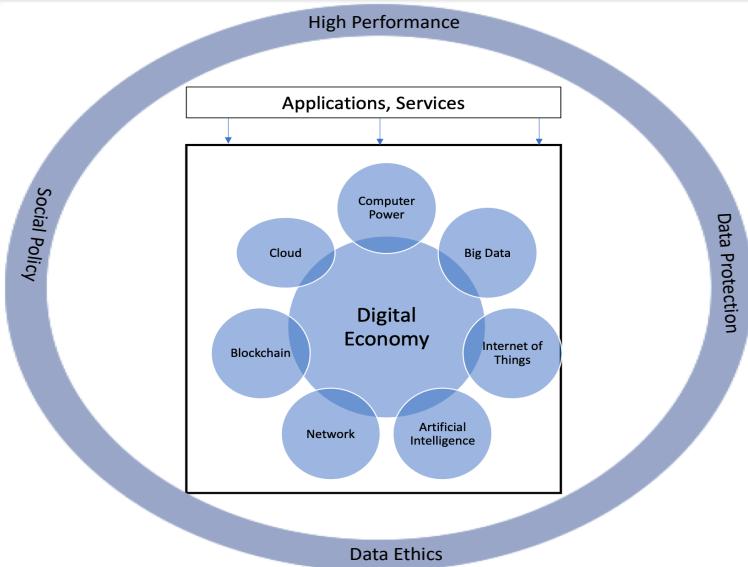
Our World  
in Data



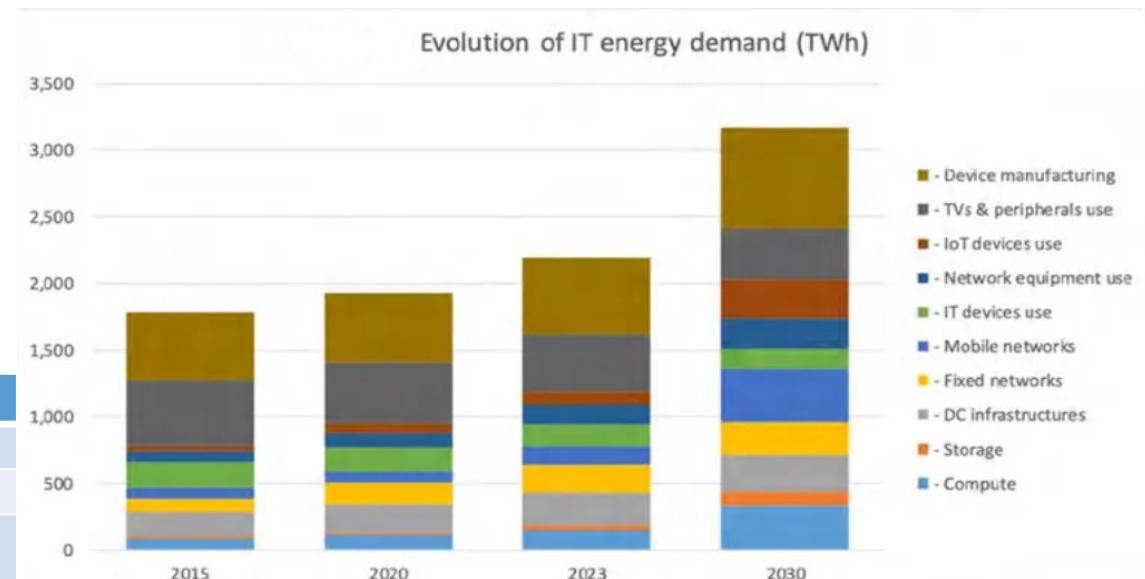
OurWorldInData.org – Research and data to make progress against the world's largest problems.  
Source: Climate Watch, the World Resources Institute (2020).

Licensed under CC-BY by the author Hannah Ritchie (2020).

# Digitalization sector



- ▶ **2030:** IT sector electricity power demand will grow by 50% by 2030
- ▶ **2040:** IT carbon footprint could reach 14%, with data centers contributing to 50% of growth



	2015	2021	Evolution
Internet users	3 billion	4,9 billion	+60%
Internet Traffic	0,6 ZB	3,4 ZB	+440%
Data center workloads	180M	650M	+260%
Data center energy use (excluding crypto)	200 TWh	220-320 TWh	+10-60%
Crypto mining energ use	4 TWh	100-140 TWh	+2300-3333%
Data transmission network energy use	220 TWh	260-340 TWh	+20-60%

<https://www.i-scoop.eu/sustainability-sustainable-development/it-sector-electricity-demand/>

# World's most valuable resource

May 6th 2017

**“Data is the new oil.”**

Clive Robert Humby  
*mathematician, entrepreneur, and Chief Data Scientist, Starcount*

**“Data is the new currency.”**

Antonio Neri, *President Hewlett Packard Enterprise*



**“Data is a commodity like gold.”**

Matt Shepherd  
*Head of Data Strategy, BBH London*

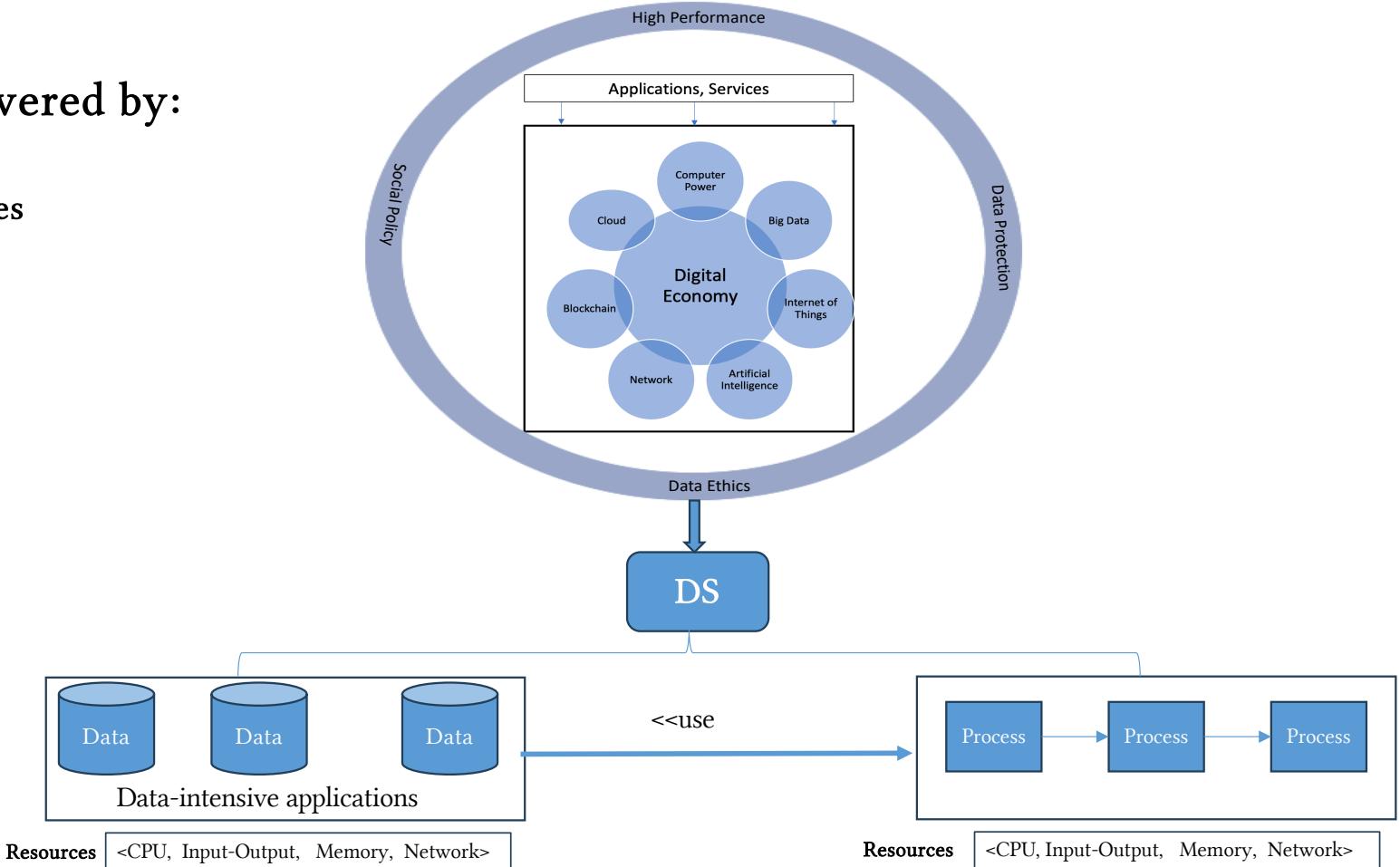
“At the heart of the digital economy and society is the explosion of insight, intelligence and information – data.

**“Data is the lifeblood of the digital economy.”**

World Economic Forum  
*A New Paradigm for Business of Data BRIEFING PAPER - JULY 2020*

# Data Science (DS) era

- ▶ Digitalisation is powered by:
  1. Data
  2. Streamlined processes

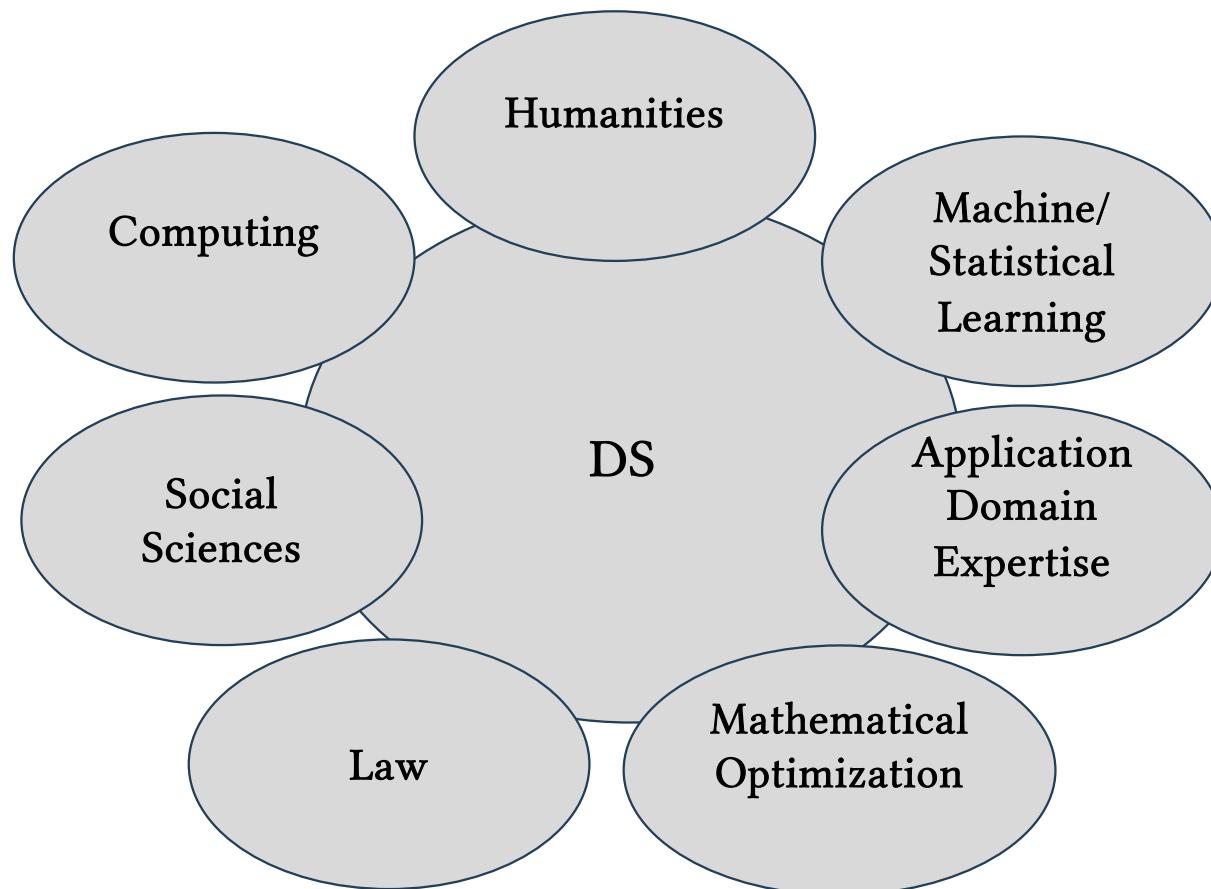


# DS: definition

- ▶ A **data-based** approach to problem solving by analysing and exploring large volumes of possibly **multi-modal data**, extracting from it knowledge and insight that is used for better **decision-making**
- ▶ It involves the **process** of collecting, preparing, managing, analysing, and explaining the data and analysis results

T. Özsü: Data Science - a systematic treatment. Communications of the ACM 66(7): 106-116 (2023)

# DS as a unifier



# DS applications

- ▶ DS is about applications (**No applications → No DS**)
  - Applications give purpose
  - Applications inform core technologies
- ▶ Almost any domain with large data sets are good candidates
- ▶ Some examples
  - Energy management (EU PLAIBDE Project)
  - Smart things (ex. cities, buildings)
  - Medicine (CHU Poitiers)
  - Recommender systems
  - Transportation
  - Fraud detection
  - Industry 4.0 (EDF)
  - Environment
  - Education & Learning (AT 41 project)
  - Fashion
  - Movie industry (CGR Cinema)
  - Finance & insurance (MAIF Foundation)
  - Food, nutrition

# DS: ecosystem

## Applications



### 4 Pillars of Data Science

#### Data Engineering

- Big data management
- Data analysis
- Data understanding
- Data preparation

#### Data Analytics

- Explore data (data mining)
- Build models & algorithms (ML)
- Visualizations & visual analytics

#### Data Protection

- Security
- Privacy

#### Data Ethics

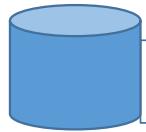
- Impact on individuals, organizations & society
- Ethical & normative concerns
- Bias in data
- Algorithmic bias
- Regulatory issues



### Social and Policy Context

Vision of GDR 3708 - MaDICS

# Data Engineering



Big data management

- Data enrichment, integration
  - Extract/Transform/Load (**joins, filters, aggregations, sorting, ...**)
  - Data lakes
- Storage and **processing** of big datasets
- Data processing platforms

Data understanding & analysis

- Data profiling
- Detection of anomalies
- Data changes

Data Preparation

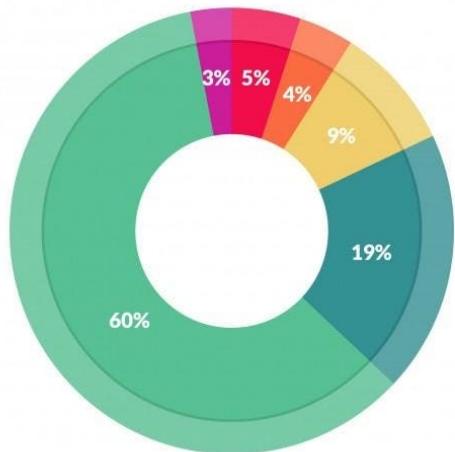
- Data acquisition/gathering
- Data cleaning
- Data provenance & lineage

# Data Engineering is important



Source: <https://www.dataquest.io/blog/advanced-data-cleaning-r-course/>

Data preparation accounts for about 80% of the work of data scientists



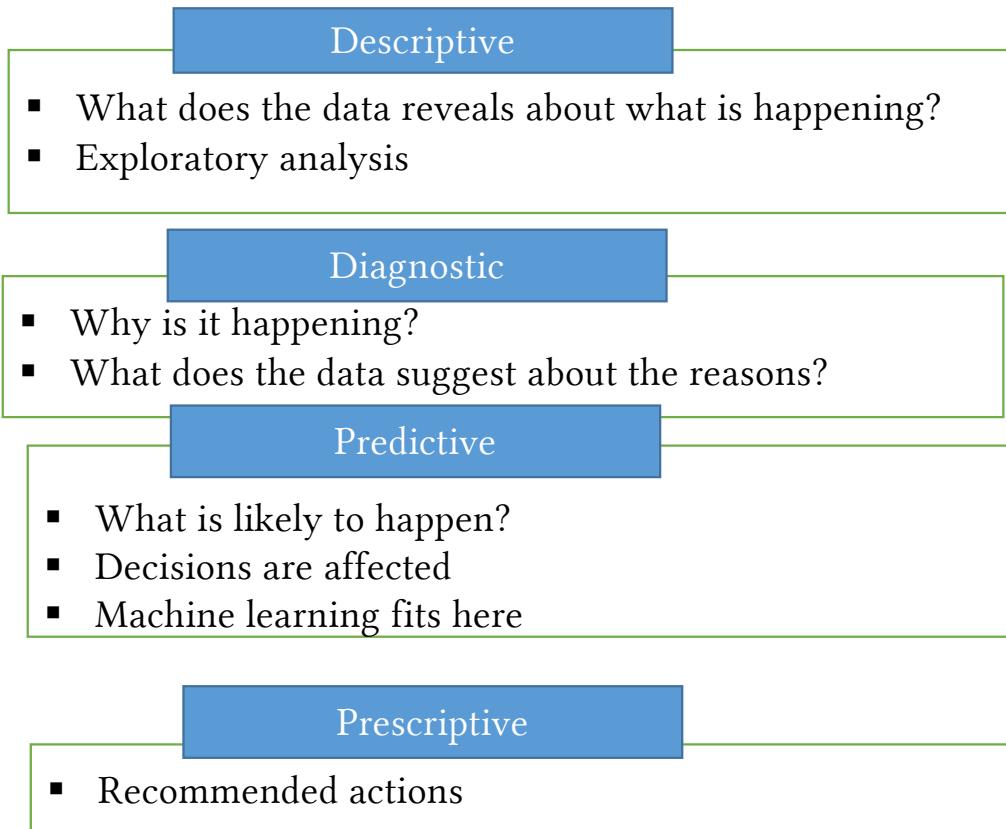
What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets; 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%

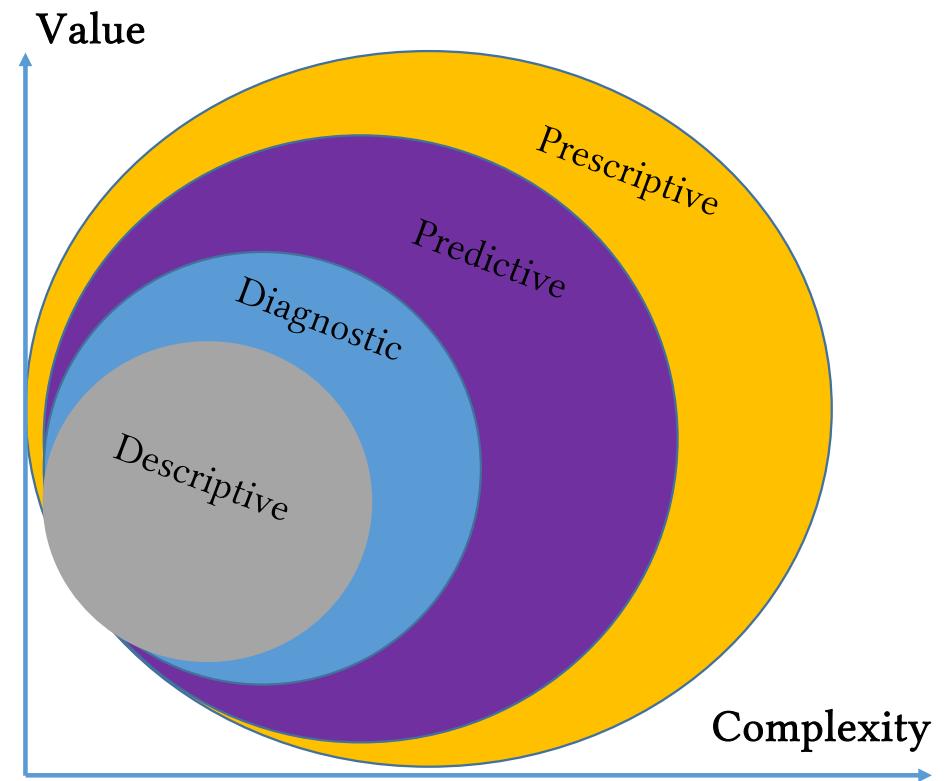
<https://www.forbes.com/sites/gilpress/2016/03/23/data-preparation-most-time-consuming-least-enjoyable-data-science-task-survey-says/>

# Data Analytics

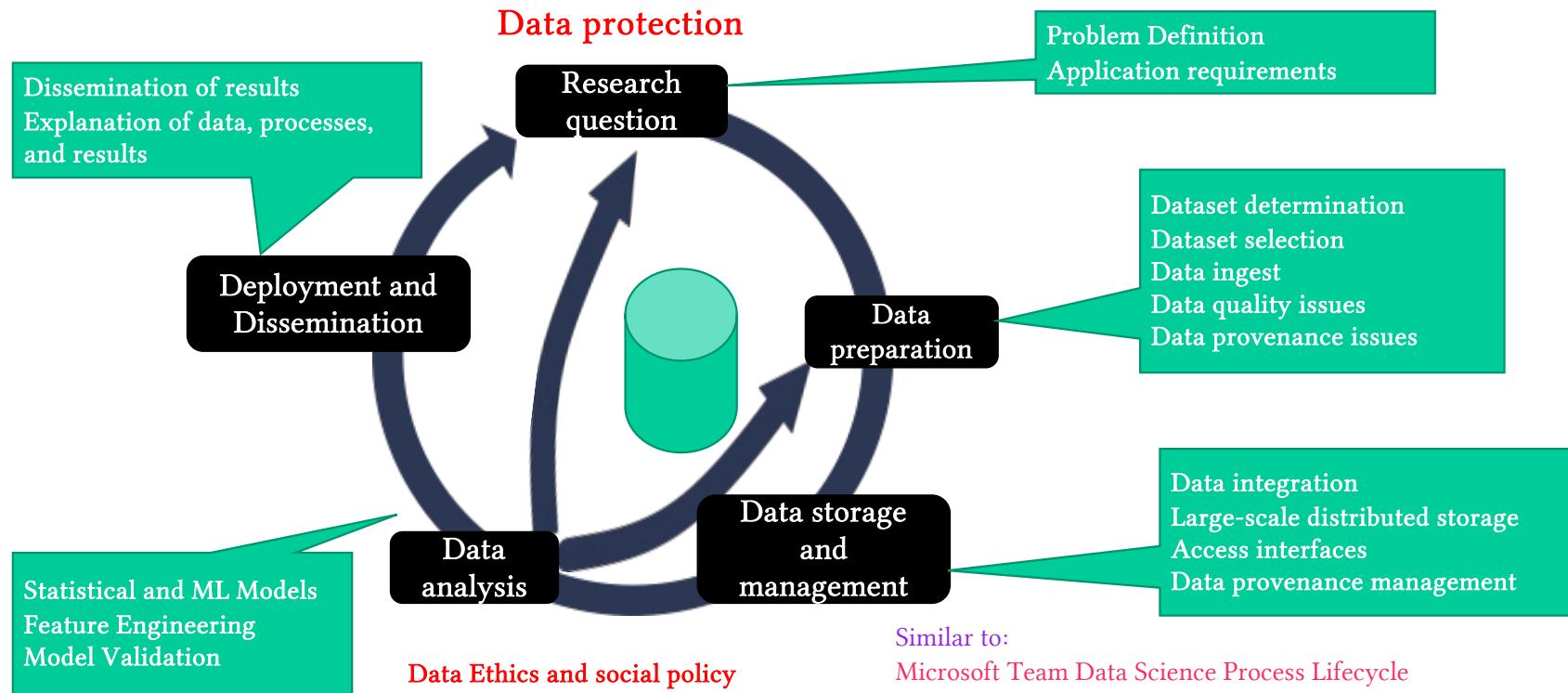
## →Drawing insights from data



<https://www.kdnuggets.com/2017/07/4-types-data-analytics.html>

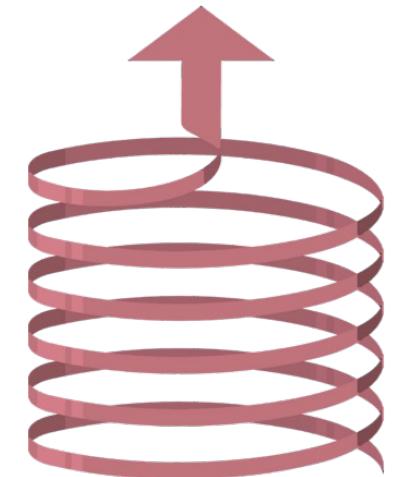
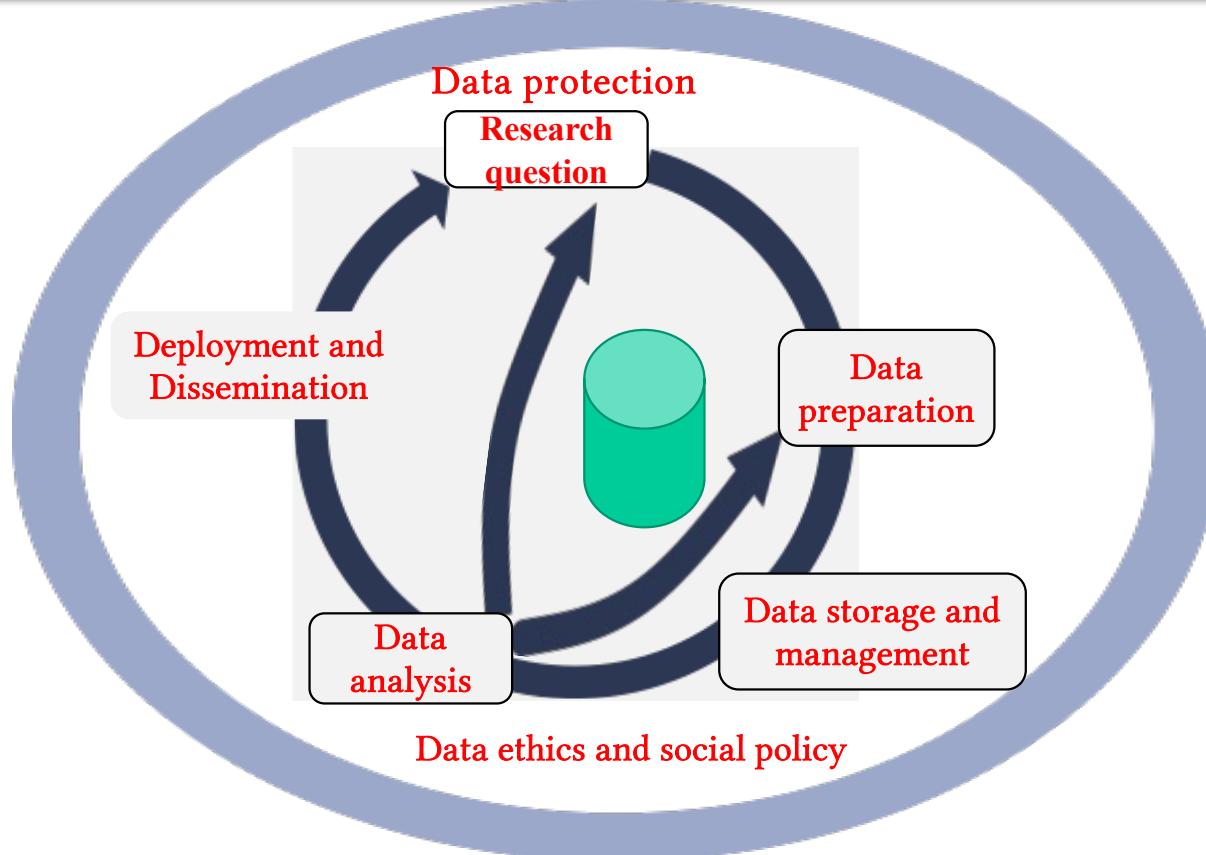


# DS life cycle (1/2)



T. Özsü: Data Science - A Systematic treatment. Communications of the ACM 66(7): 106-116 (2023)

## DS life cycle (2/2)

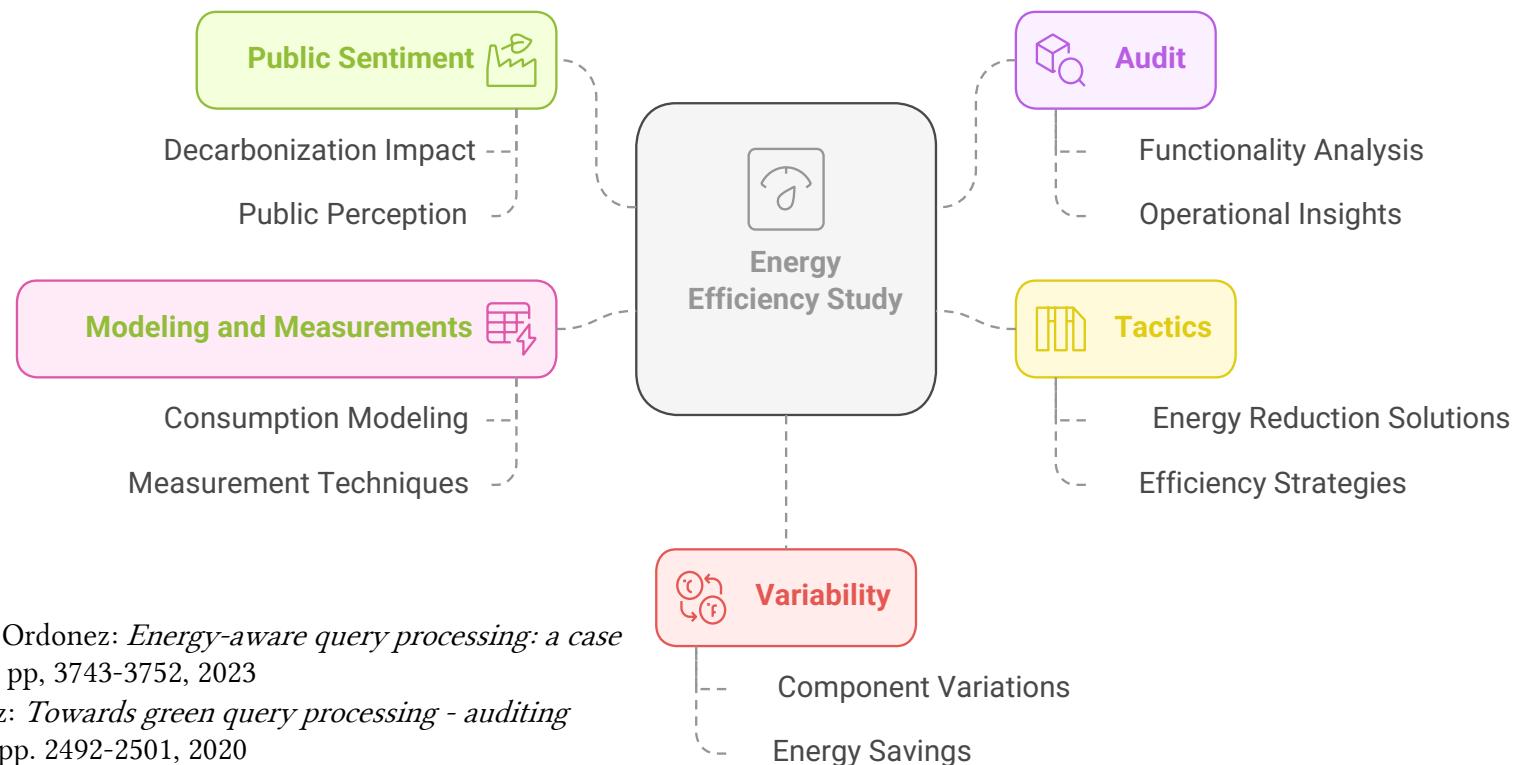


- ▶ The development of DS applications often requires significant time and resources (ex., energy)
- ▶ Like any other sector, DS needs to be examined from an EE perspective

# Generic framework for studying EE

## STAMP-V

- Defined based on our experience in studying EE (e.g., Nearly Zero Energy Building, Databases)
- STAMP-V is designed to assess EE of any object (e.g., house, car, airplane, server)



1. L. Bellatreche, F. Djellali, W. Macyna, C. Ordóñez: *Energy-aware query processing: a case study on join reordering*. **IEEE Big Data**, pp. 3743-3752, 2023
2. S. P. Dembele, L. Bellatreche, C. Ordóñez: *Towards green query processing - auditing power before deploying*. **IEEE Big Data**, pp. 2492-2501, 2020

# STAMP-V (Digitalisation)

## ► Public Sentiment

### ■ Two controversial questions

1. How can digitalisation be utilized to mitigate the causes and impact of CC?
2. Does the digitalisation sector has a negative impact on the environment and how can it be reduced?



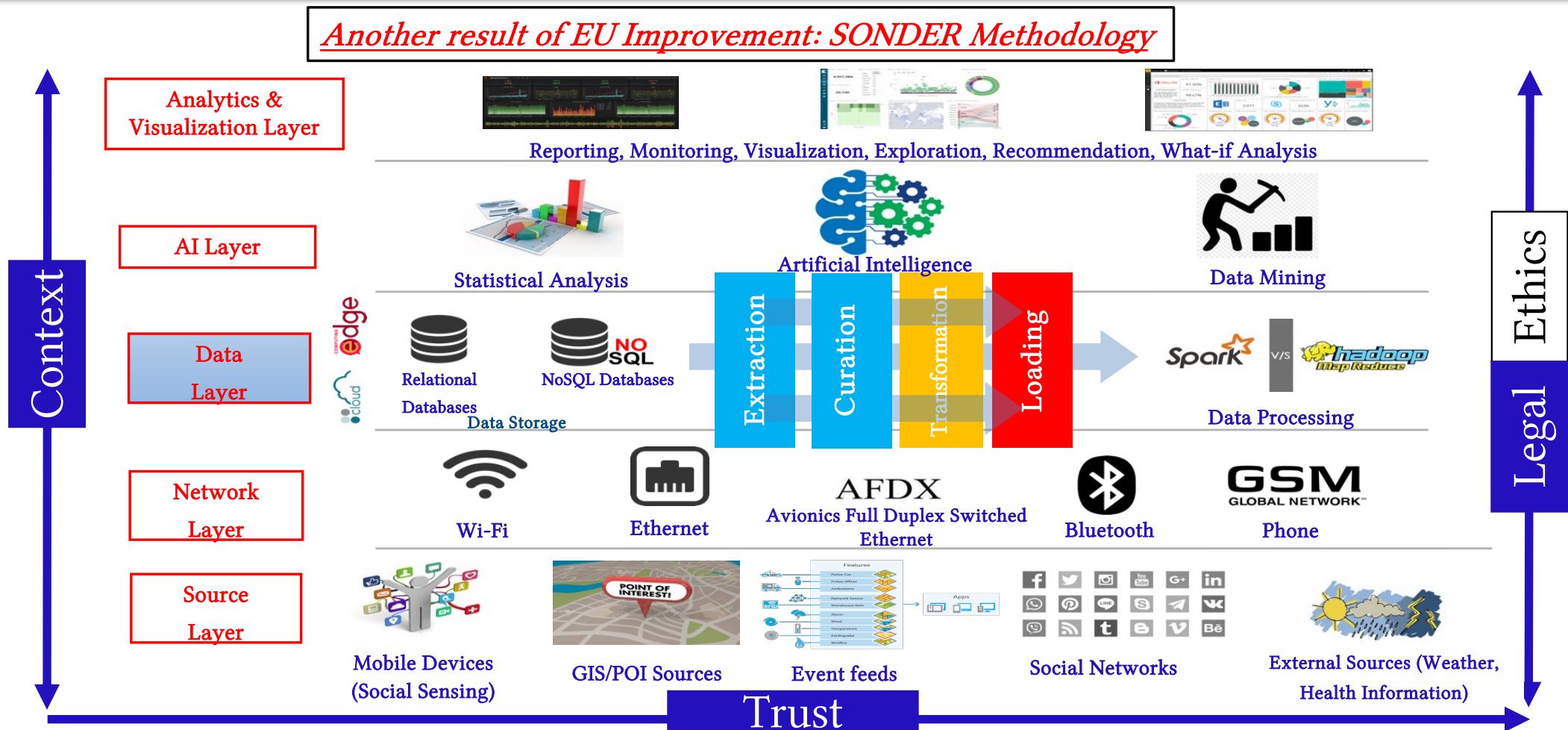
## ► Success stories: DS to the rescue

1. BCOOLER of Google's DeepMind team enhanced the efficiency of cooling systems in data centers by 12.7% ([using Reinforcement Learning](#))
2. Zero Energy Public Buildings: EU IMPROVEMENT
  - Data-driven solutions (ML/DL techniques)
    - o Energy Consumption and Price Prediction



1. J. Tobajas, F. Garcia-Torres, P. Roncero-Sánchez, J. Vázquez, L. Bellatreche, E. Nieto, *Resilience-oriented schedule of microgrids with hybrid energy storage system using model predictive control*, *Applied Energy Journal*, 118092, 306, Part B, 2022
2. F. Chauvet, L. Bellatreche, C. A. S. Silva: *AI approaches for electricity price forecasting in stable/unstable markets: EU improvement project*. *IEEE Big Data*, pp. 4473-4482 , 2022

# STAMP-V (Digitalisation)



# STAMP-V (Digitalisation)

“Data is the new oil.”

Clive Robert Humby  
*mathematician, entrepreneur, and Chief Data Scientist, Starcount*



”Data isn’t the new oil, it’s the new CO2

Martin Tisné, Vice President, Luminate Strategic Initiatives

→AI harming the environment 😞

- Cloud data centers (data storage and processing) consumed approximately 205 terawatt-hours (TWh) of electricity in 2020 ( $\approx 1\%$  of the world’s total electricity consumption)
- Heating, ventilation, and air conditioning are responsible for a significant percentage of global CO2 emissions
- Workflow applications are energy intensive: training, building and deploying models

Model	Number of features	CO2 equivalent emissions
Gopher	280B	322 tons
BLOOM	176B	26 tons
GPT3	175B	502 tons
GPT	175B	70 tons

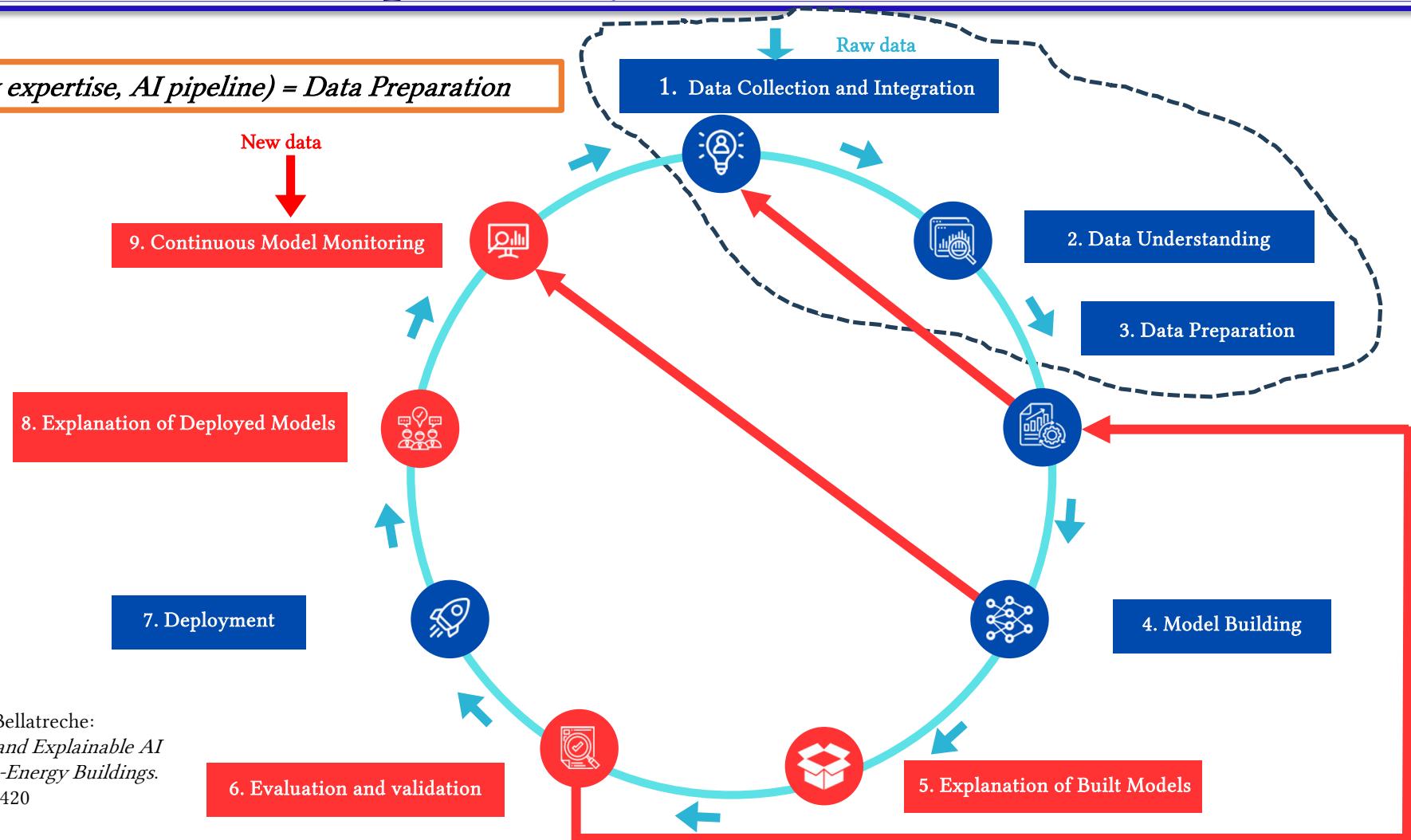
Source: Luccioni et al., 2023

→What can we do in this critical situation?

- Each scientist has to be energy-aware

# Scope of my initiatives

*Projection (my expertise, AI pipeline) = Data Preparation*

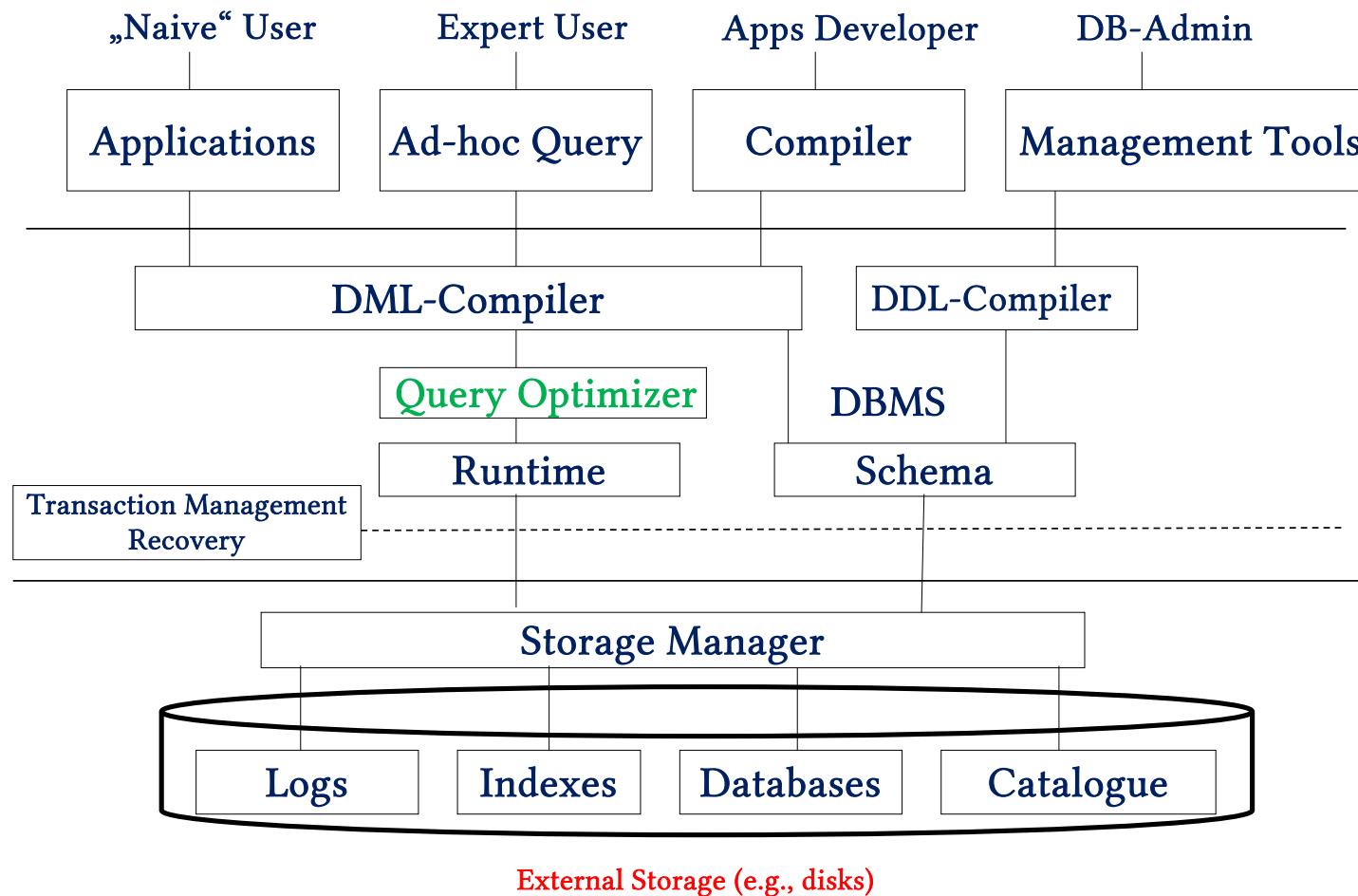


# Action plans

1. QP have become the backbone for data preparation and data exploration phases, enabling the efficient use of ML/DL techniques
  2. Think big, start small: a good initiative to design green QP
  3. Reuse of our physical design experience
    - a. Quantification of energy consumption of a QP
    - b. Building green QP
    - c. *Green physical design*



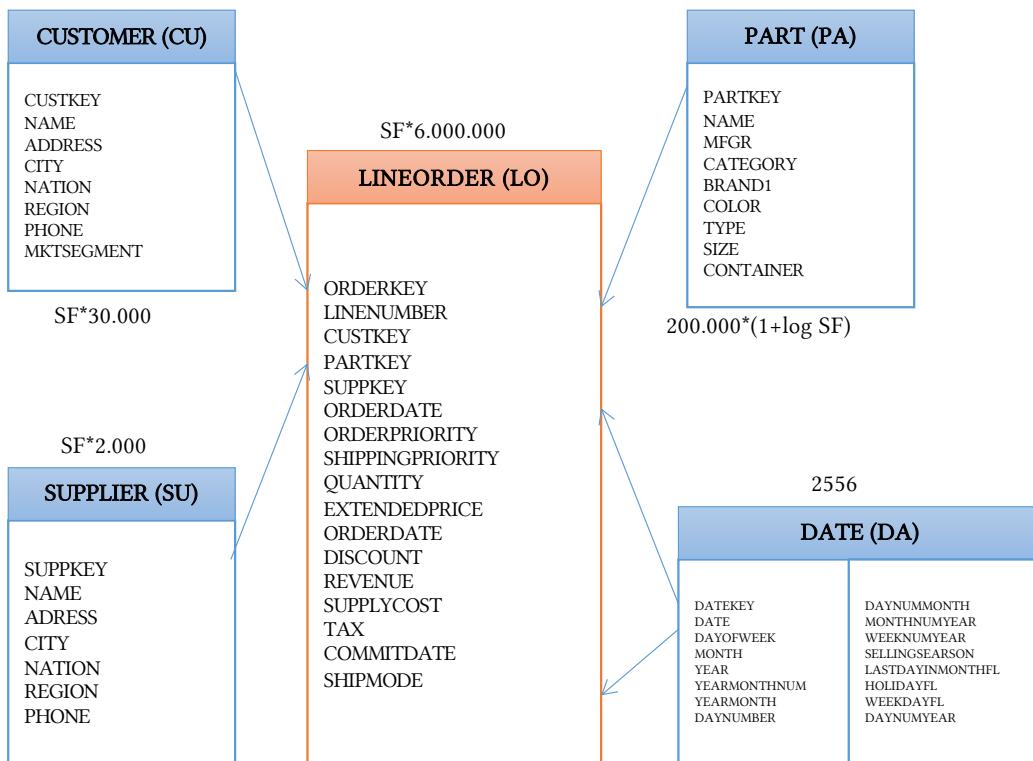
# STAMP-V (QP)



# STAMP-V (QP)

## Analytical SQL queries

### Star Schema Benchmark (SSB)



### Query example

```

SELECT c_nation, s_nation, d_year, sum(lo_revenue)
FROM customer CU, lineorder LO, supplier SU, dates DA
WHERE lo_custkey = cu_custkey          (J1)
AND lo_suppkey = su_suppkey          (J2)
AND lo_orderdate = da_datekey       (J3)
AND cu_region = 'AMERICA'           (S1)
AND su_region = 'AMERICA'           (S2)
AND da_year ≥ 1993 AND da_year ≤ 1998 (S3)
GROUP BY cu_nation, su_nation, da_year
ORDER BY da_year asc, lo_revenue desc;

```

↓

Data-intensive (joins, selection, sorting, aggregations)

# STAMP-V (QP)

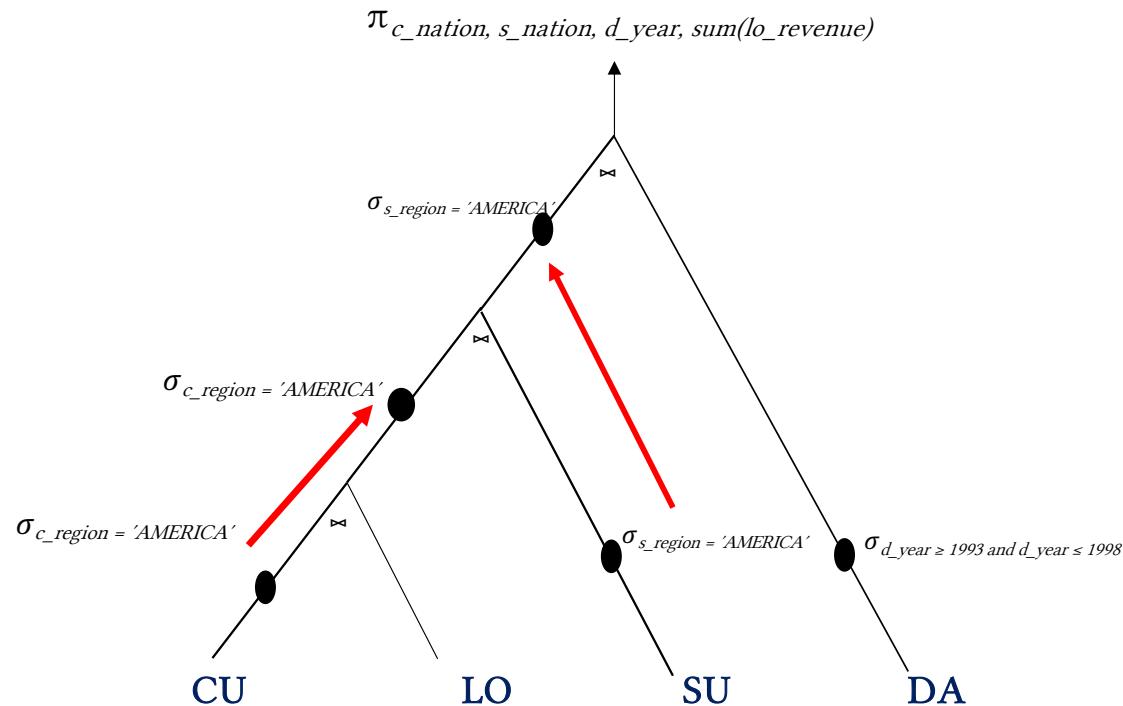
```
import pandas as pd  
  
# Sample DataFrame 1  
df1 = pd.DataFrame({  
    'ID': [1, 2, 3, 4], 'Name': ['Alice', 'Bob', 'Charlie', 'Diana'] })  
  
# Sample DataFrame 2  
df2 = pd.DataFrame({'ID': [2, 3, 4, 5], 'Department': ['HR', 'Engineering', 'Sales', 'Finance'], 'Salary': [50000, 70000, 60000, 80000]})  
  
# Merge type: 'inner', 'left', 'right', 'outer'  
merge_type = 'inner'  
  
# Merge on 'ID'  
merged_df = pd.merge(df1, df2, on='ID', how=merge_type)  
  
# 🔎 Filter: only employees with salary > 55000  
filtered_df = merged_df[merged_df['Salary'] > 55000]  
  
# Sort by Salary descending  
sorted_df = filtered_df.sort_values(by='Salary', ascending=False)  
  
# Display final result  
print("Merged + Filtered + Sorted DataFrame:\n")
```

## AI Job: Pandas world

% Join, filter, sorting are present

# STAMP-V (QP)

- ▶ Finding an optimal plan is a hard problem
  - Variety of query trees (cascading of selections, join ordering, ...)



# STAMP-V (QP)

## A query plan

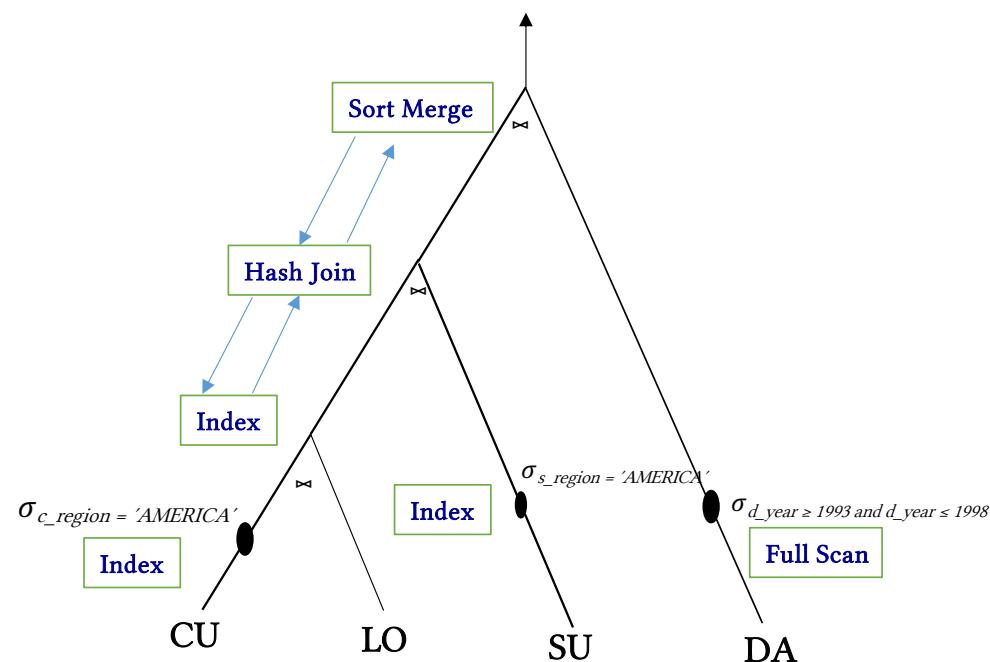
- Query tree in which each node has is annotated by its **implementation algorithm**
- Multiple algorithm per node

### I. Unary operations

- Selection
  - Sequential scan
  - Index

### II. Binary operations

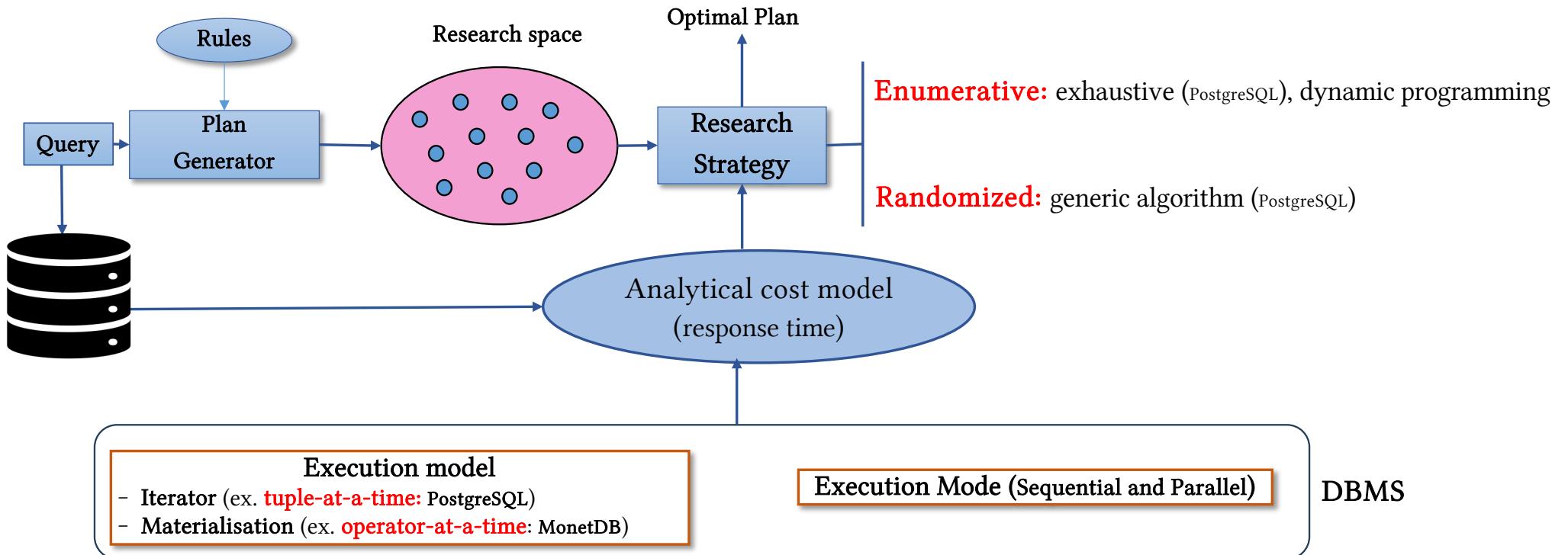
- Join
  - Nested-loop
  - Indexed nested-loop
  - Sort-Merge
  - Hash
  - ...
- Choice based on cost estimate

$$\pi_{c\_nation, s\_nation, d\_year, \text{sum}(lo\_revenue)}$$


→ Problem of finding the *best plan*  
▪ NP-hard

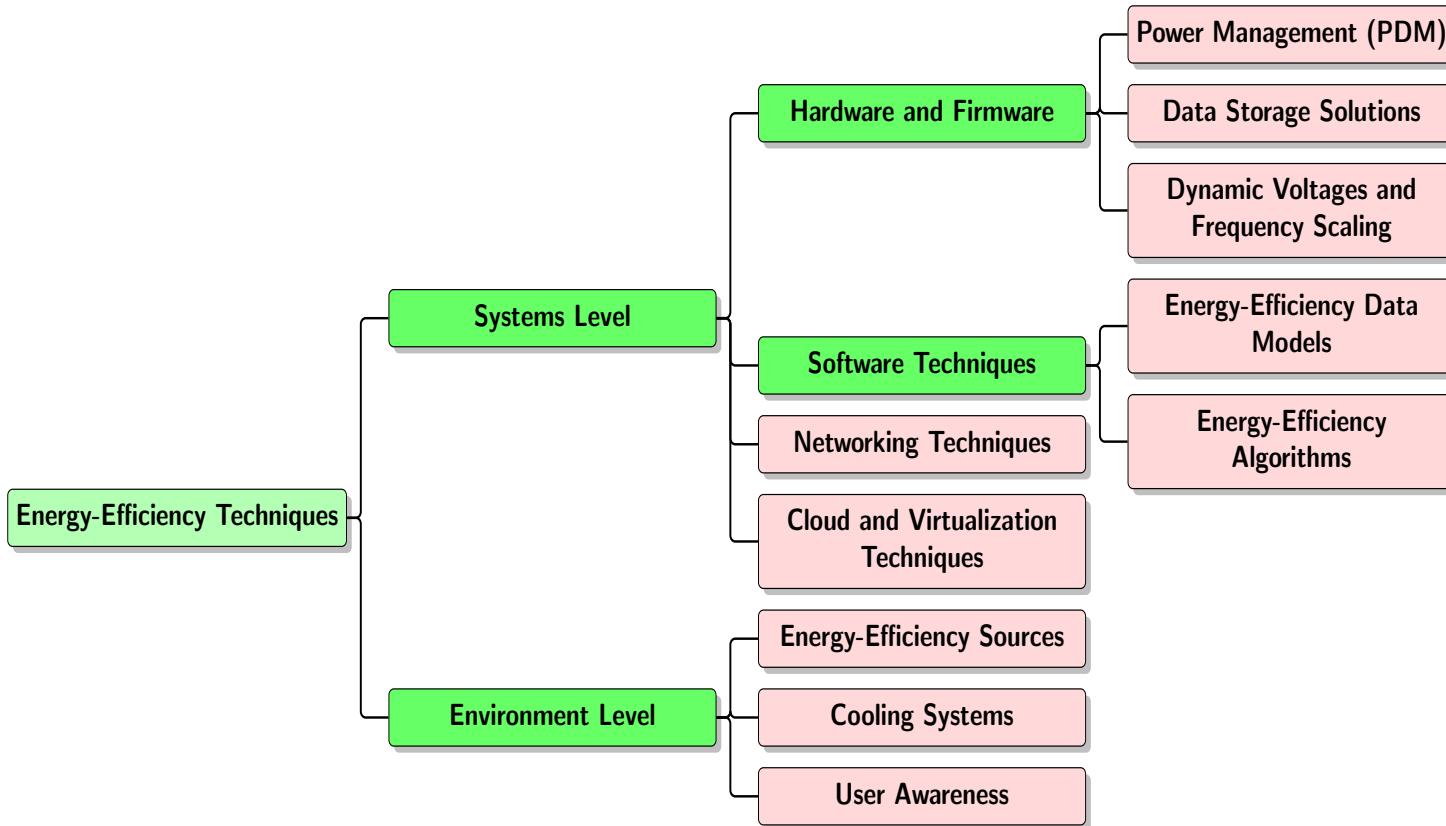
# STAMP-V (QP)

- ▶ Cost-based approach for selecting optimal plan



→ Use this audit to integrate energy consumption into query processor

# STAMP-V (QP)



# Hardware Tactics

- ▶ **Dynamic Power Management**
    - Dynamic component deactivation
    - Dynamic performance scaling
  - ▶ **Hardware accelerators**
    - Graphical Processing Unit (GPU)
    - Field-Programmable Gate Array (FPGA)
    - Tensor Processing Units (TPU)
  - ▶ ...

L. Bellatreche, W. Macyna, and C. Ordóñez, *Green Analytics*, Tutorial, IEEE Big Data Conference 2023

# Hardware Tactics

- ▶ Observation: power states can be dynamically adjusted to meet current performance requirements

## Dynamic component deactivation

disabling relevant hardware component(s) when they are idle

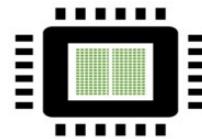
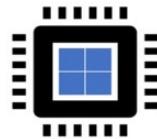
## Dynamic performance scaling

the clock frequency of certain components, such as the CPU, can be decreased or increased

## Hardware Tactics

- ▶ **Dynamic Voltage and Frequency Scaling (DVFS) – one of the most popular techniques in DynPS**
  - ▶ GOAL: It allows a system to adjust the frequency and supply voltage to particular components within the infrastructure computation
  - ▶ DVFS on a multi-core CPU
    - System with only one global voltage for all cores (global DVFS) is energy-inefficient.
    - Global DVFS and per-core DVFS architectures with multiple Voltage Frequency Islands (VFIs) have been proposed.
    - The cores in an island share the same voltage and frequency, but different islands can be executed at various voltages and frequencies

# Hardware Tactics



CPU	GPU
Central Processing Unit	Graphical Processing Unit
4-32 cores	100s or 1000s of cores
Low latency	High throughput
Good for serial processing	Good for parallel processing
Quickly process tasks that require interactivity	Break jobs into separate tasks to process simultaneously
Traditional programming are written for CPU sequential execution	Requires additional software to convert CPU functions for parallel execution
Less power to operate but low performance	More power to operate but high performance

The energy advantage of the GPU varies between approximately 2 and 4 times, depending on the workload.

# Software Tactics

- ▶ SaaS is the most popular economic model used by data science in the cloud.
  - ▶ Energy should be included in this model by adopting the "polluter pays" principle.
  - ▶ This approach will encourage cloud providers to be more environmentally responsible.

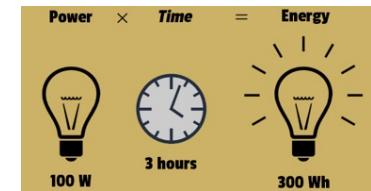


*In 2020, the global SaaS market was valued at \$120.77 billion and is expected to reach \$462.94 billion by 2028.*

- ▶ Monitoring and feedback loops
  - ▶ Efficient algorithms and data structures
  - ▶ Workload categorization and prediction
  - ▶ Query scheduling
  - ▶ Buffer management
  - ▶ Virtualisation and containers

# ~~STAMP~~-V: Modeling & Measurement

- ▶ **Energy** (E, joules, watts-hour): the ability to do work.
- ▶ **Power** (P, watts): the rate of energy (E) per a unit of time (T)
- ▶ **Baseline Power:** power consumption when the machine is idle
- ▶ **Active Power:** power consumption due to the execution of the workload
- ▶ **Peak power:** represents the maximum power
- ▶ **Average power:** average power consumed during execution
- ▶ **Power Usage Effectiveness (PUE):** Total Power/IT Power



# ~~STAMP~~-V: Modeling & Measurement

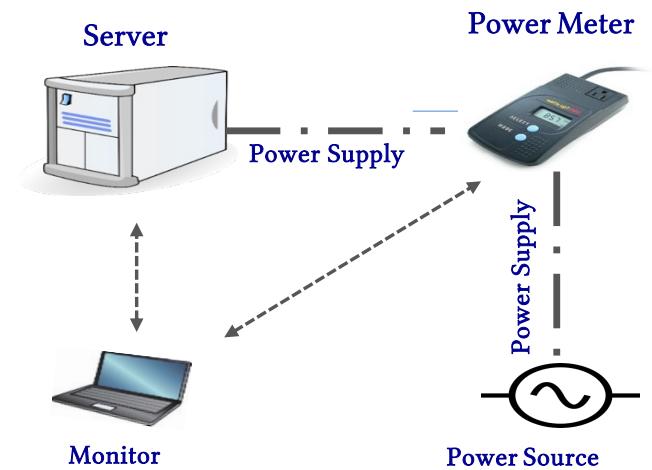
## 1. Real measurements

- Database as a blackbox
- Use of monitor the power consumption of electrical device (ex. yocto-watt)

## 2. Estimated Measurements: CodeCarbon



- An open-source tool to estimate carbon emissions from computing
- Developed by Mila, BCG GAMMA, and Haverford College
- Supports Python workflows
- Works with CPU, GPU, and cloud providers
- Monitors (CPU/GPU usage, execution time, location-based carbon intensity)
- Calculates (power consumed, carbon dioxide equivalent (CO<sub>2</sub>eq) emitted)
- Outputs (emission logs (.csv), visualization dashboard)



# ~~STAMP-V~~: Modeling & Measurement



Reusing expertise from response time cost models

- ▶ Additive energy measurement models
- ▶ Energy consumption of an algorithm A ( $E(A)$ )
  - $E(A) = E_{\text{cpu}}(A) + E_{\text{memory}}(A)$
- ▶ Energy consumption of a server E(S)
  - $E(S) = E_{\text{cpu}}(S) + E_{\text{memory}}(S) + E_{\text{I/O}}(S)$

S. Roy, A. Rudra, and A. Verma, “An energy complexity model for algorithms,” in Proc. 4th Conf. ITCS, 2013, pp. 283–304.

# ~~STAMP~~ V: Modeling & Measurement

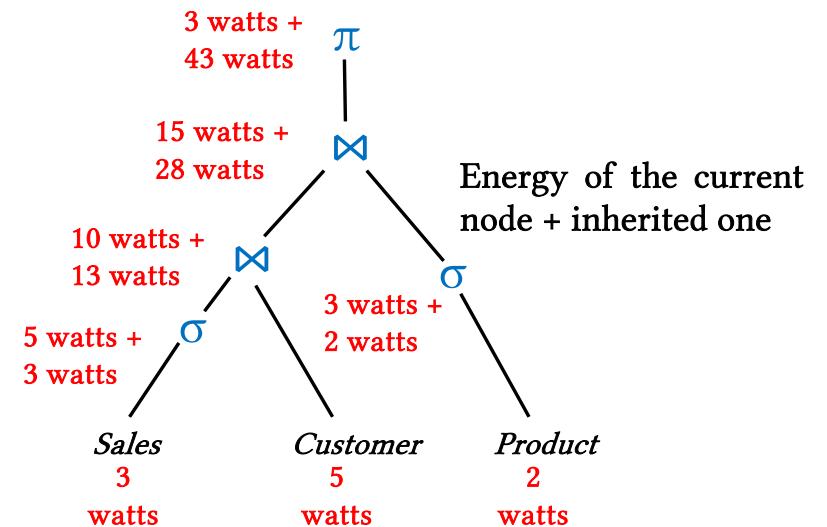
## ► Energy consumption cost model

- CPU, Inputs-Outputs (IO), network

## ► Principle

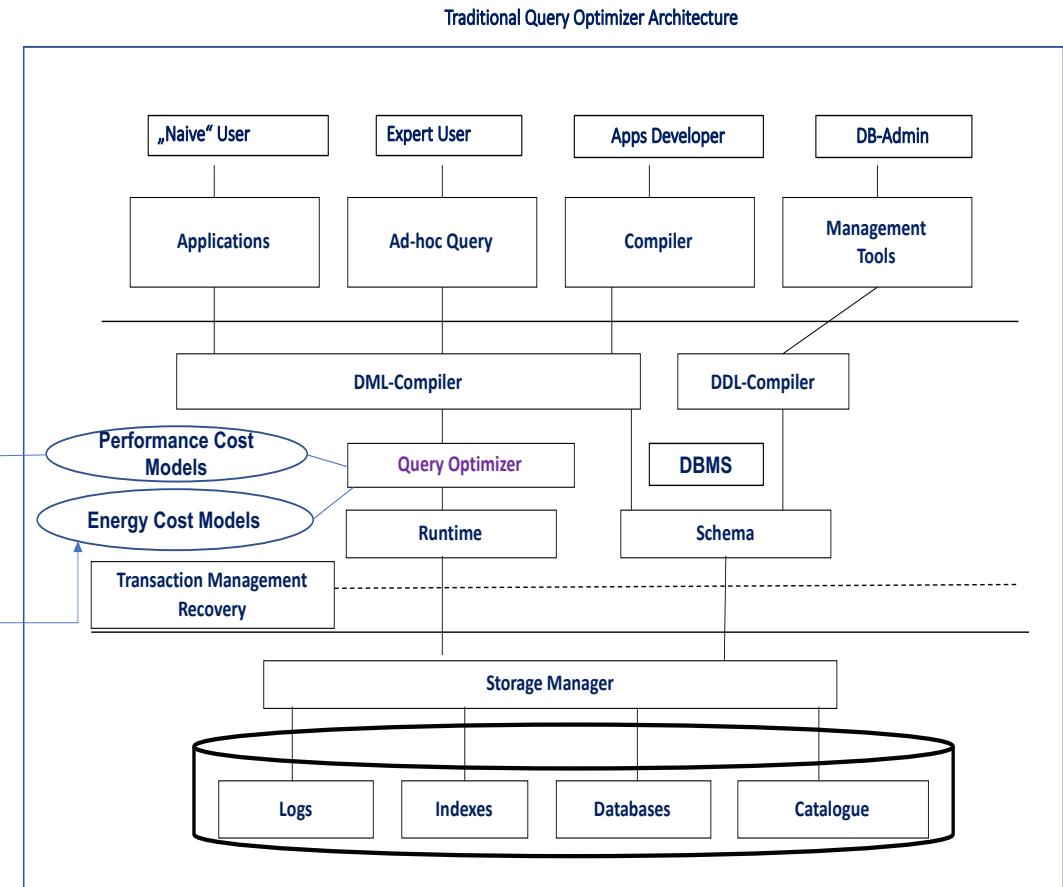
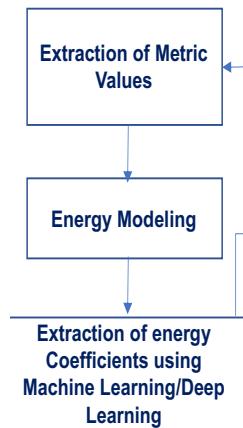
- Understand algorithms used by each SQL operator (join, sort, ...)
  - estimate available main memory buffers
  - estimate the size of inputs, intermediate results
- Composition of cost of operators:
  - Aggregation of resource consumption

Total power: 46 watts



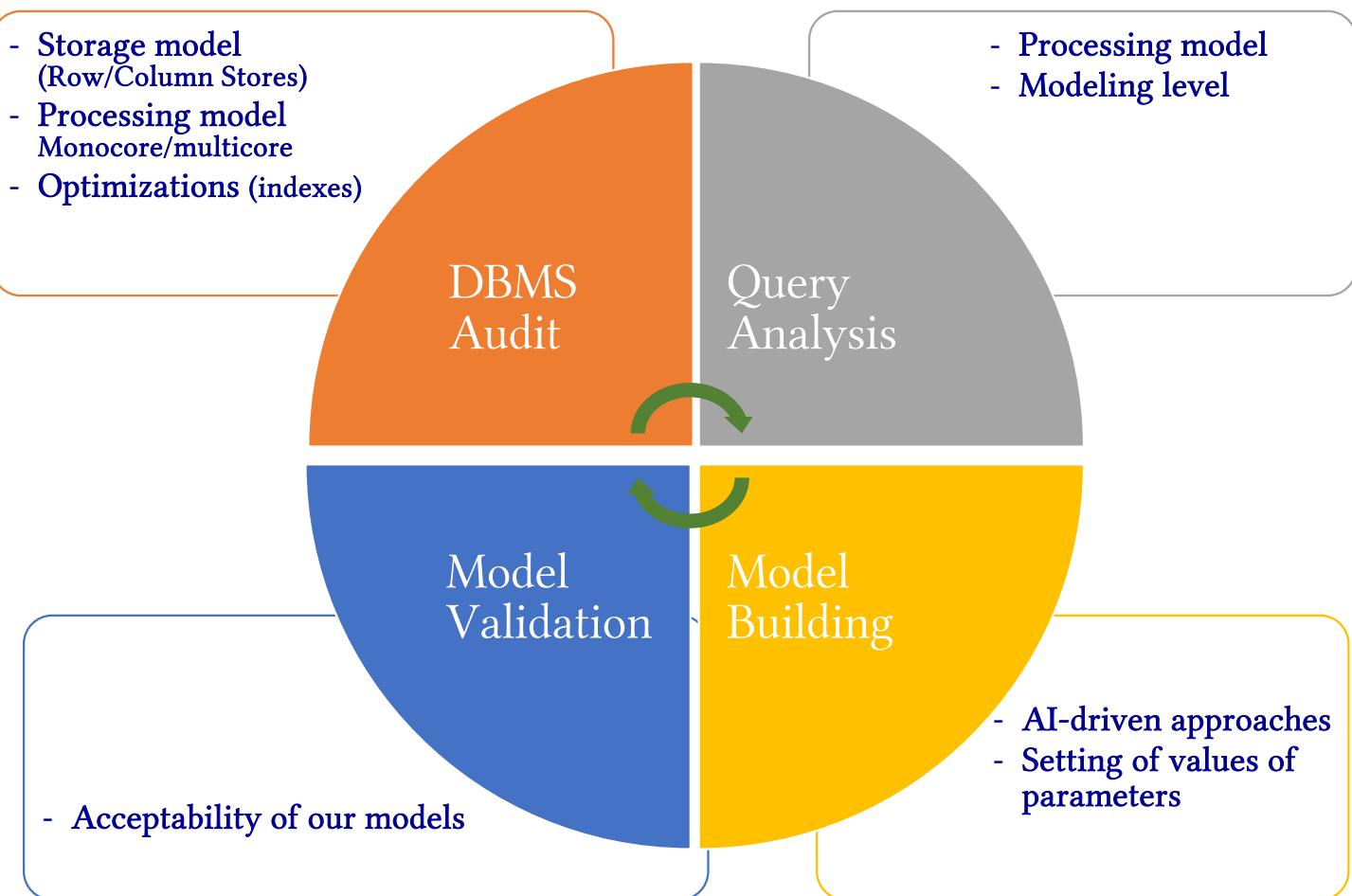
# ~~STAMP~~-V: Modeling & Measurement

- ▶ How to estimate consumed energy of each operator?
  - Vision 1: On the top of DBMS vision
  - AI-Driven Approaches to Estimation



# ~~STAMP~~-V: Modeling & Measurement

→ Our modelling approach



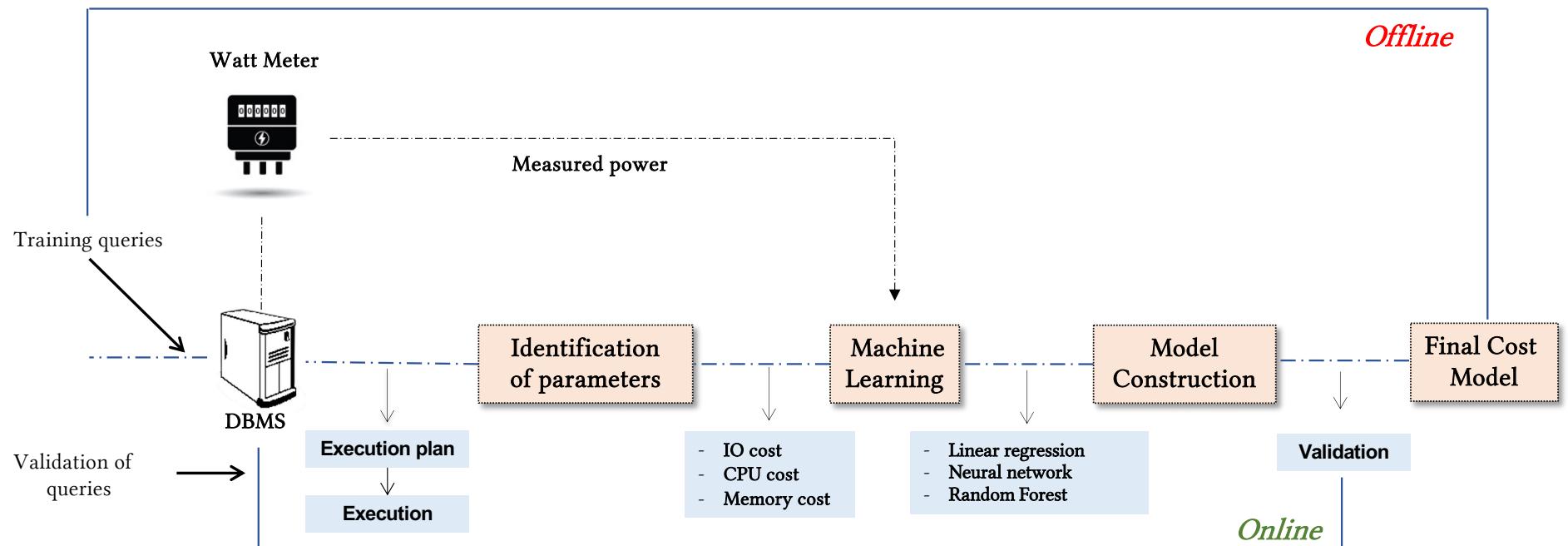
# ~~STAMP~~ V: Modeling & Measurement

► Inputs

- DB schema
- Training queries
- DBMS
- Platform

► Used DBMSs

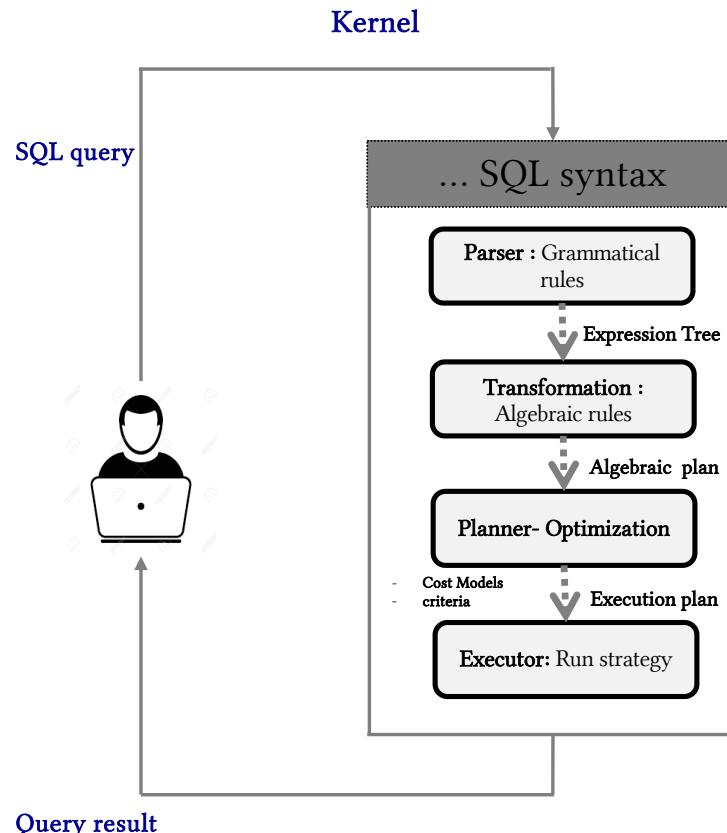
- PostgreSQL : Explain, discard all, ...
- MonetDB : Explain, Tracer, plan, PAPI, PCM
- Watt Meter



# ~~STAMP-V~~: Modeling & Measurement

## Deployment 1: PostgreSQL

 PostgreSQL: Widely used in Research and Teaching (OLTP Applications)



- Major audit information

1. Row-oriented storage layout

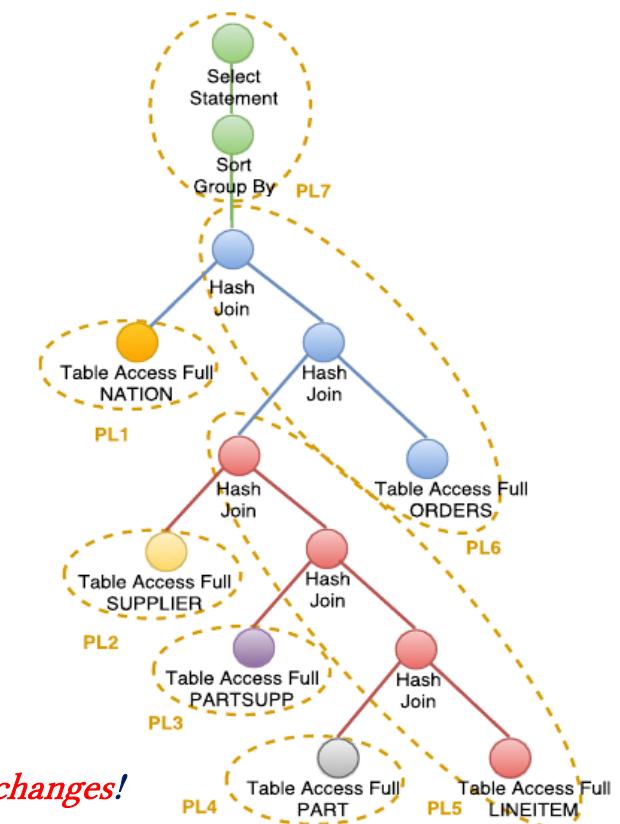
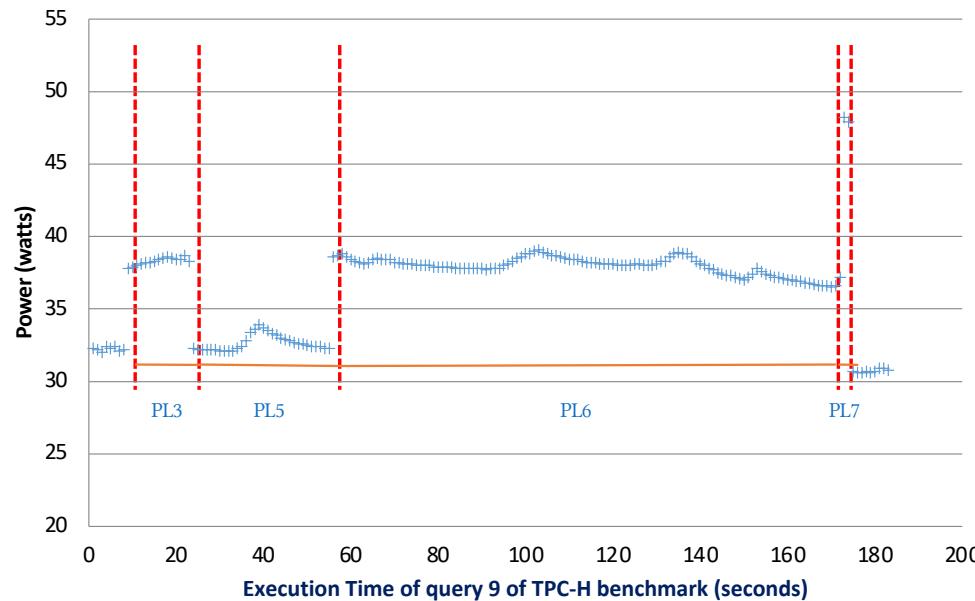
CID	CName	Gender	Age
1	Simon	M	36

2. Only SQL queries
3. Parallel query processing (multicore)
4. Modelling level: pipeline
5. Other parameters (indexes, selectivity of predicates, buffer politics, size of tables, pages of table, join implementations, ...)

→ Energy Cost = CPU + IO

# STAMP-V: Modeling & Measurement

## Deployment 1: PostgreSQL



- When a query execution goes from one pipeline to another, its energy consumption changes!

A. Roukh, L. Bellatreche: Eco-processing of OLAP complex queries. DaWaK, pp. 229-242, 2015

# ~~STAMP~~-V: Modeling & Measurement

- Machine Learning

$$\text{Power}(Q_i) = \frac{\sum_{j=1}^p \text{Power}(PL_j) * \text{Time}(PL_j)}{\text{Time}(Q_i)}$$

$$\text{Power}(PL_j) = W_{cpu} * \sum_{u=1}^n C_{cpu_u} + W_{mem} * \sum_{u=1}^n C_{mem_u} + W_{dio} * \sum_{u=1}^n C_{dio_u}$$

Parameters	Meanings
$C_{cpu_u}$	Number of instructions executed by CPU
$C_{mem_u}$	Number of read/write operations accessing memory
$C_{dio_u}$	Number of read/write operations accessing disk
n	Number of operators in the pipeline

→ all provided by PostgreSQL

→ ML Technique to calculate  $\beta_i$ : *multiple polynomial regression (degree=2)*

$$\begin{aligned} P(PL_j^i) = & \beta_1 * C_{cpu} + \beta_2 * C_{mem} + \beta_3 * C_{dio} + \beta_4 * C_{cpu} * C_{mem} + \beta_5 * C_{mem} * C_{dio} \\ & + \beta_6 * C_{cpu} * C_{dio} + \beta_7 * C_{cpu}^2 + \beta_8 * C_{mem}^2 + \beta_9 * C_{dio}^2 + \beta_0 + \varepsilon \end{aligned}$$

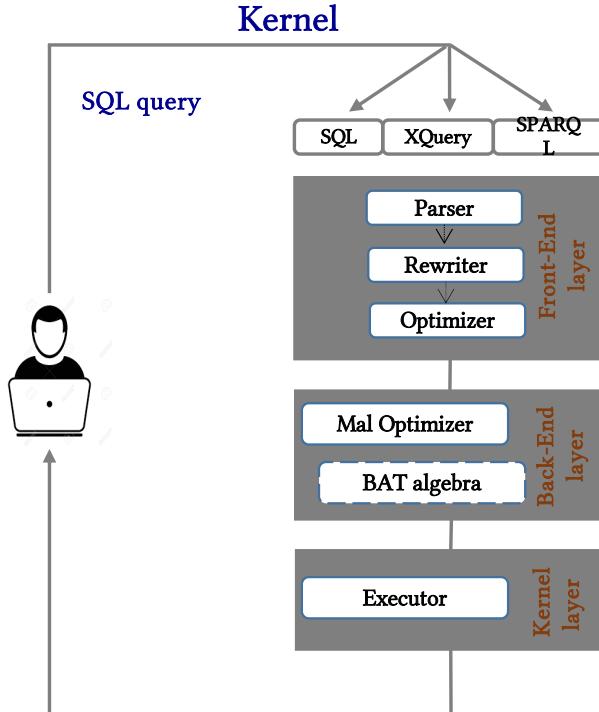
- $\beta_i$ : Regression coefficients
- $\varepsilon$ : Measurement error

# STAMP V: Modeling & Measurement

## Deployment 2: MonetDB



### Analytical Queries



Query result

- Major audit information

1. Column-oriented storage layout

CID	CName	Gender	Age
1	Ahmed	M	36

2. Compression

3. Multi query languages

4. Parallel query processing (multicore)

5. Modelling level: operation

6. ...

→ Energy Cost = CPU + IO

# ~~STAMP-V~~: Modeling & Measurement

- Machine Learning

$$\text{Power}(Q_i) = \frac{\sum_{j=1}^n \text{Power}(OP_j^i) * \text{Time}(OP_j^i)}{\text{Time}(Q_i)}$$

$$\text{Power}(OP_j^i) = W_{cpu} * C_{cpu_j} + W_{mem} * C_{mem_j} + W_{dio} * C_{dio_j}$$

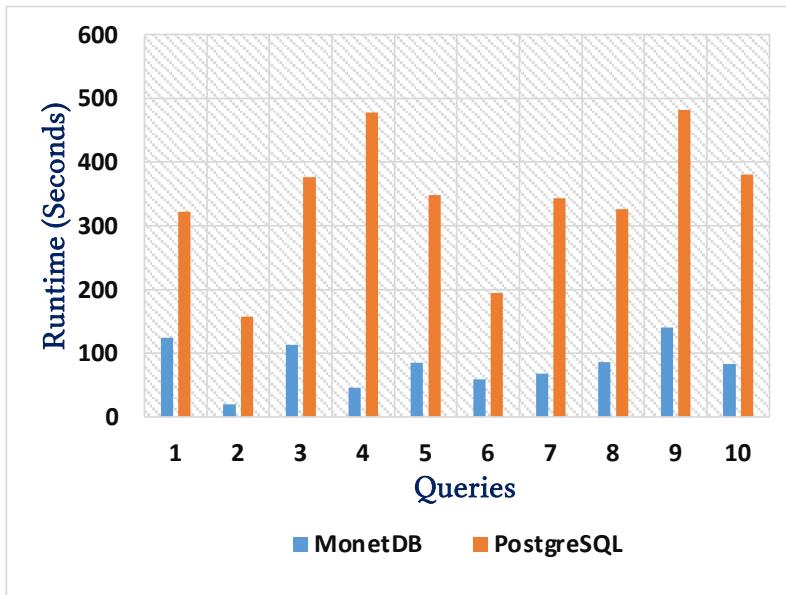
Parameters	Meanings	
$C_{cpu_u}$	Number of instructions executed by CPU	→ Provided by Processor Counter Monitor
$C_{mem_u}$	Number of read/write operations accessing memory	
$C_{dio_u}$	Number of read/write operations accessing disk	→ provided by MonetDB

→ Machine Learning Technique to calculate  $\beta_1$ : *Multiple Polynomial Regression (degree = 3)*

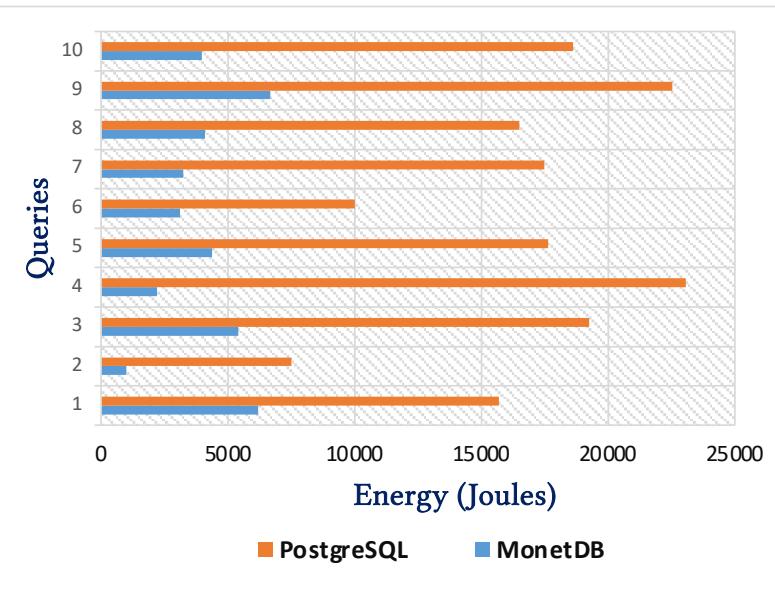
$$\begin{aligned} P(OP_j^i) = & \beta_1 * C_{cpu} + \beta_2 * C_{mem} + \beta_3 * C_{dio} + \beta_4 * C_{cpu} * C_{mem} + \beta_5 * C_{mem} * C_{dio} \\ & + \beta_6 * C_{cpu} * C_{dio} + \beta_7 * C_{cpu} * C_{mem} * C_{dio} + \beta_8 * C_{cpu}^2 + \beta_9 * C_{mem}^2 \\ & + \beta_{10} * C_{dio}^2 + \beta_{11} * C_{cpu}^2 * C_{mem} + \beta_{12} * C_{mem}^2 * C_{cpu} + \dots \\ & + \beta_{17} * C_{cpu}^3 + \beta_{18} * C_{mem}^3 + \beta_{19} * C_{dio}^3 + \beta_0 + \varepsilon \end{aligned}$$

- $\beta_i$  : Regression coefficients
- $\varepsilon$ : Measurement error

# ~~STAMP~~-V: Modeling & Measurement



*Execution time of queries*  
[30 GB data set of TPC-H benchmark]

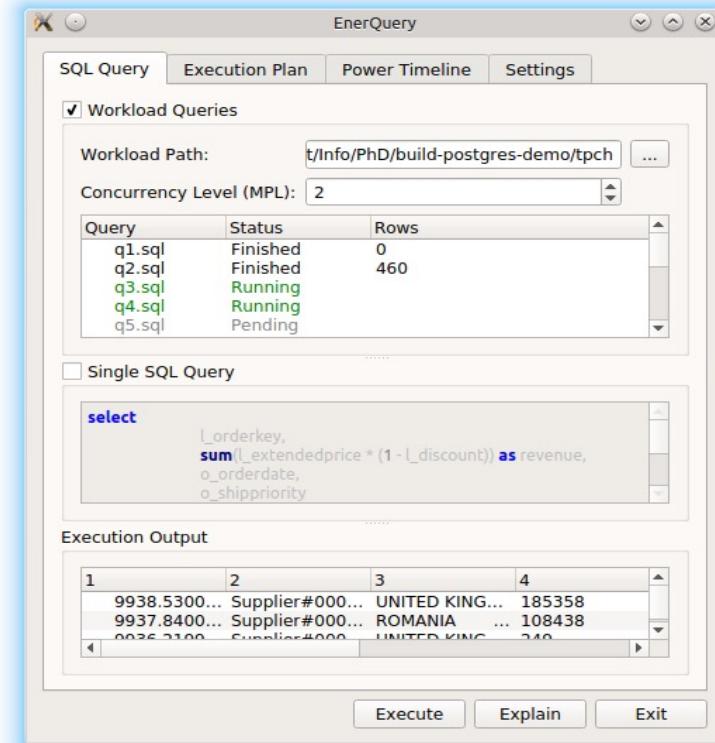
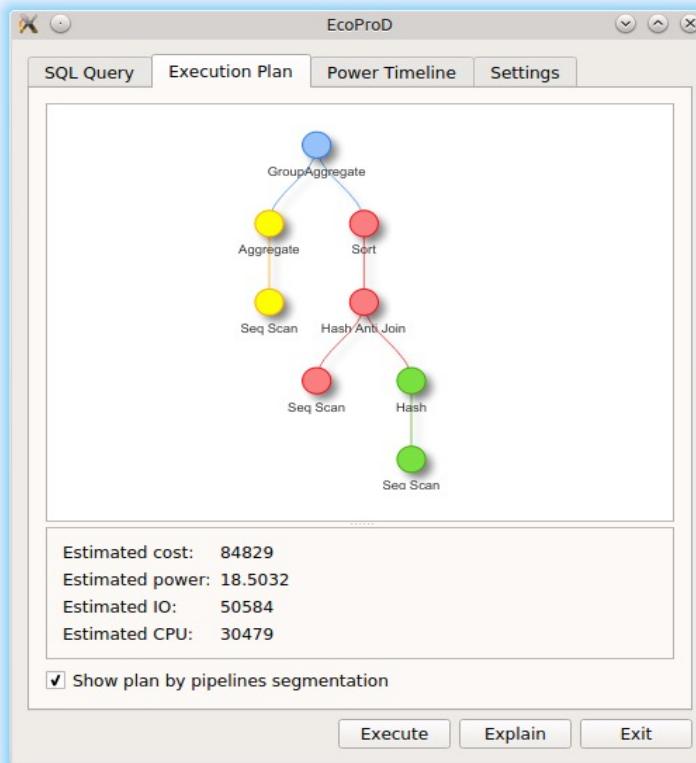


*Energy consumed of queries queries*  
[30 GB data set of TPC-H benchmark]

**MonetDB consumes less energy than PostgreSQL**

# ~~STAMP~~-V: Modeling & Measurement

- ▶ EnerQuery: <https://youtu.be/cWK5rf4MBNQ>
- ▶ Available at: <https://github.com/lias-laboratory/enerquery>



# STAMP-V : Variability

Work	S	A	T	M	V
Tsirogiannis et al.	No	Data Processor	Software & Hardware	Real Measurement	Black Box
Bouhatous et al.	No	Data processor	Software & Hardware	Analytical cost models and ML	Black Box (DBMS) + Variation of Processor frequency
Dembele et al.	No	Data processors	Software	Analytical cost models and ML	Black Box
Roukh et al.	No	Data processor	Software	Analytical cost models and ML	Black Box
Höpfner et al.	No	Data processor	Software	Real Measurement	Join and sorting implementation
Guo et al.	No	Data processor	Software	Analytical cost models and ML	Black Box

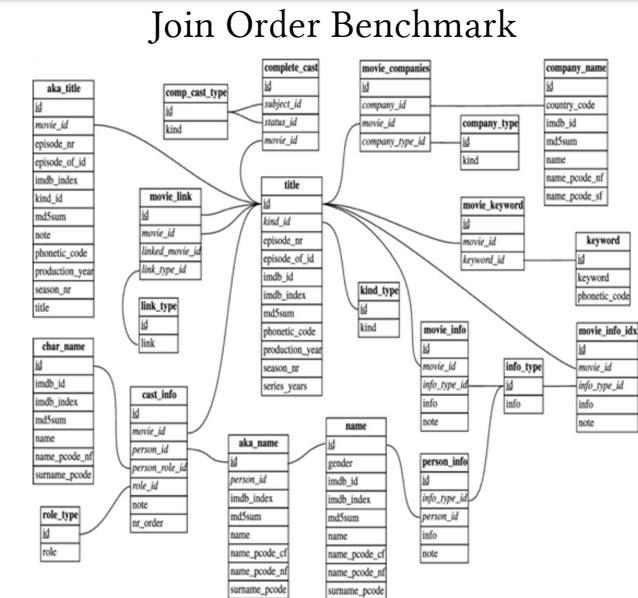
# STAMP-V : Variability

- ▶ The number of possible join trees grows rapidly as the number of join relations increases
- ▶ If a query involves 21 tables, then employing left deep tree, the number of potential join orders is 21! ( $5.1090942e+19$ )
- ▶ NP-hard problem
- ▶ A QP is responsible for identifying the optimal join strategy and join implementation

## Diversity of join ordering algorithms

**1980 – 2020**

1. Exact solutions
2. Randomized Solutions (e.g., genetic algorithm, simulated annealing)
3. Mixed integer linear programming
4. Heuristics (e.g., minimum selectivity, Monte Carlo Tree Search)



**2020 - Now**

1. Data-driven algorithms (e.g., RTOS)

# ~~STAMP~~-V: Variability

## → Variation of join ordering algorithm

### Query Q (Star Join Benchmark)

```
SELECT MIN(n.name) AS member_in_charnamed_movie
FROM cast_info AS ci, company_name AS cn, keyword AS k, movie_companies AS
mc, movie_keyword AS mk, name AS n, title AS t
WHERE k.keyword ='character-name-in-title'
AND n.id = ci.person_id AND ci.movie_id = t.id AND t.id = mk.movie_id
AND mk.keyword_id = k.id AND t.id = mc.movie_id AND mc.company_id = cn.id
AND ci.movie_id = mc.movie_id AND ci.movie_id = mk.movie_id
AND mc.movie_id = mk.movie_id AND n.name like '%S%';
```

### Join order 1: PostgreSQL ((k mk) t ci) n mc cn))

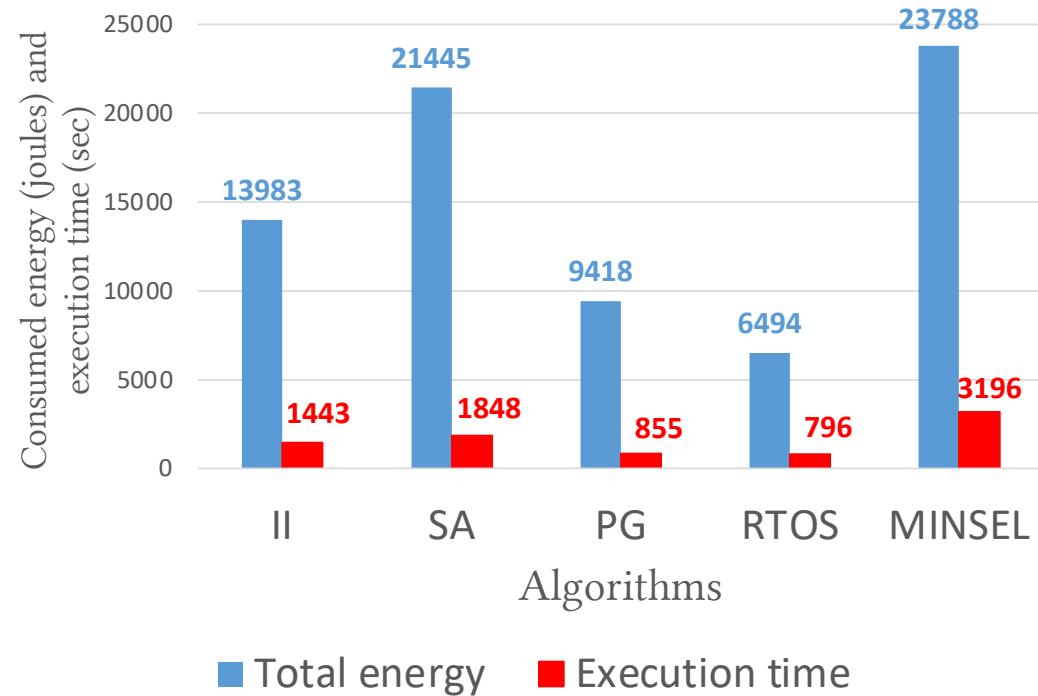
- Execution time: 17s 606ms
- Energy consumption: 0.074J

### Join order 2: (((k mk) mc) cn) ci) n) t))

- Execution time: 2s 104ms
- Energy consumption: 0.044 J

1. Join Order 2 outperforms PostgreSQL's join order in terms of query response time and energy consumption
2. The re-execution of Q using PostgreSQL gives the same costs; its query optimizer did not learn from its previous errors

# ~~STAMP~~-V: Variability



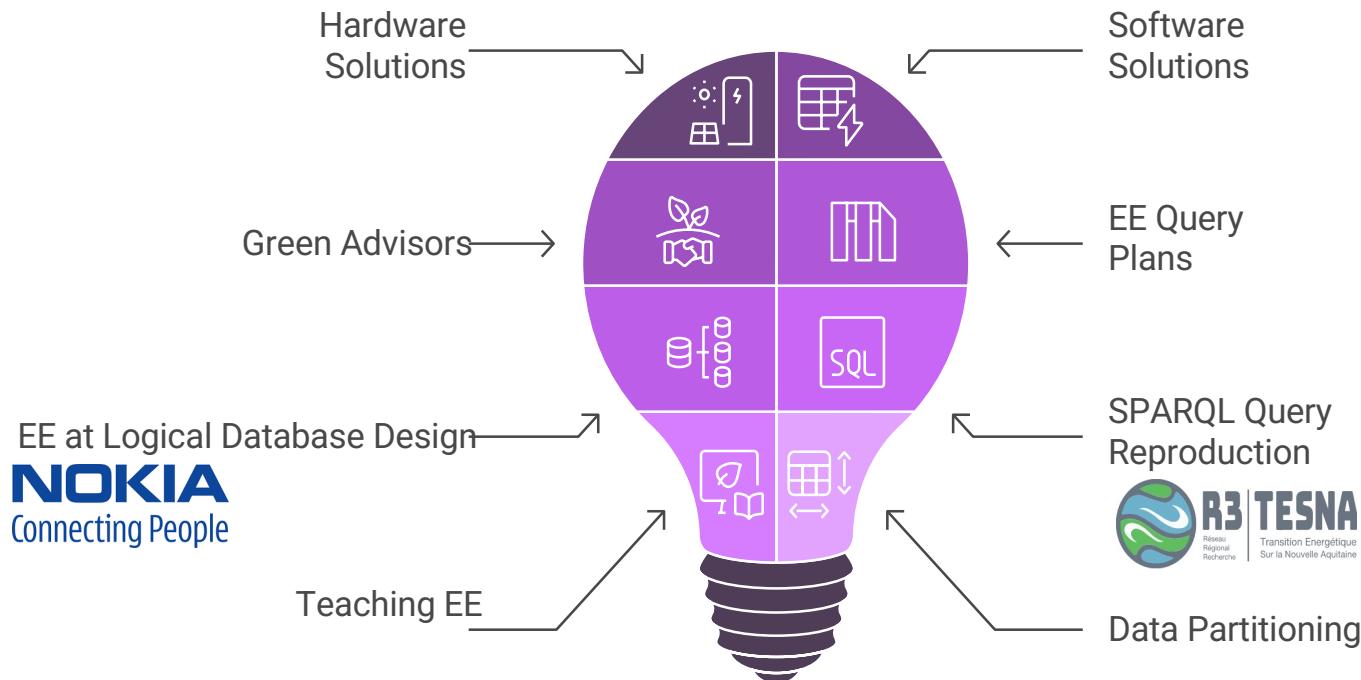
II: Iterative Improvement  
SA: Simulated Annealing  
PG: PostgreSQL  
RTOS: ML approach  
MinSel: Minimum Selectivity

Time and energy for JOB workload

# Summary

- Energy awareness is the **responsibility** of every scientist
  - EE in digitalization and AI is a **niche area**
  - DS's **positive** and **negative** impacts on energy consumption
  - STAMP-V: A comprehensive **framework** for studying EE
    1. Sentiment Analysis,
    2. Audit,
    3. Tactics,
    4. Modeling & Measurements
    5. Variability
  - Application of STAMP-V to QP
  - Studying EE requires a deep understanding of hardware and software

# Perspectives



# Thank to:



Dr. Amine ROUKH  
(2017)



Dr. Simon Pierre DEMBELE  
(2021)



Dr. Issam GHABRI  
(2022)



Dr. Ahcène BOUKORCA  
(2016)



Dr. Abdelkader OUARED  
(2019)



Dr. Selma BOUARAR  
(2016)

*p*  
*h*  
*d*

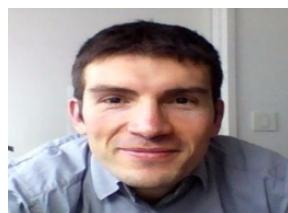
*s*  
*t*  
*u*  
*d*  
*e*  
*n*  
*t*  
*s*



Prof. Carlos ORDONEZ  
University of Houston  
USA



Prof. Sadok BEN YAHIA  
University of Southern Denmark



Prof. Stéphane JEAN,  
Poitiers University  
France

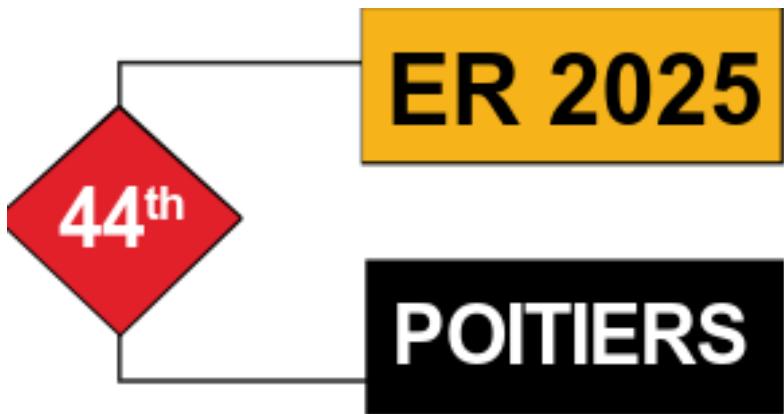


Dr. Wojciech MACYNA  
Wroclaw University  
Poland



# ER @ Poitiers

---



- Main Conference
- ER Workshops
- ER Industrial Track
- ER Tutorials
- ER Posters and Demo
- ER Forum
- ER Doctorial Consortium

44th International Conference on Conceptual Modeling

20-23 October 2025

Poitiers / Futuroscope, France

<https://er2025.ensma.fr/>

# Questions? Suggestions? Criticisms?

---



[bellatreche@ensma.fr](mailto:bellatreche@ensma.fr)

[www.lias-lab.fr/members/bellatreche](http://www.lias-lab.fr/members/bellatreche)