# Data Integration and Large Scale Analysis

## 01- Introduction and Overview

Lucas Iacono. PhD. - 2024

TU Graz
SCIENCE
PASSION
TECHNOLOGY

ISDS

PUBLIC DOMAIN

# How to contact me?

\_ \_ \_

- liacono@know-center.at (for now…)

- lucasiacono

- SAL Building – 2nd Floor – Office 02 067

- https://lucasiacono.github.io/ws2024_2025_dia.html

# Agenda

- Organization
- Motivation and Goals
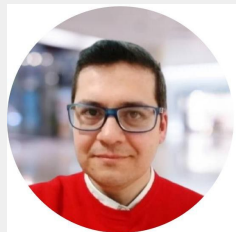- Learning Units and Projects

# About Me

# About Me

———

- **Degree in Electrical and Electronics Engineering** – University of Mendoza, Argentina (2007)
- **Doctor in Engineering (PhD)** – University of Mendoza, Argentina (2015)
- **Data Engineer** – Fiat Petronas PSG16 Race Team (2015)
- **Associate Professor** –
  - Introduction to Technology – National University of Cuyo (2015 – 2019) – Argentina
  - Industrial Robotics – Computer Engineering – University of Mendoza (2008 – 2019) – Argentina
  - Postdoc – IoT Devices and UAVs applied to frost damage mitigation – Argentine Research Council (CONICET) – 2017
- **Know Center Research GmbH**
  - Senior Researcher – HAI Area (2019 – 2023)
  - Research Area Manager – Data Management for AI Area (2023 – Now)
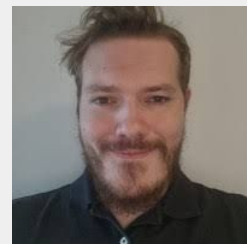
# Data Management for AI @Know Center

———

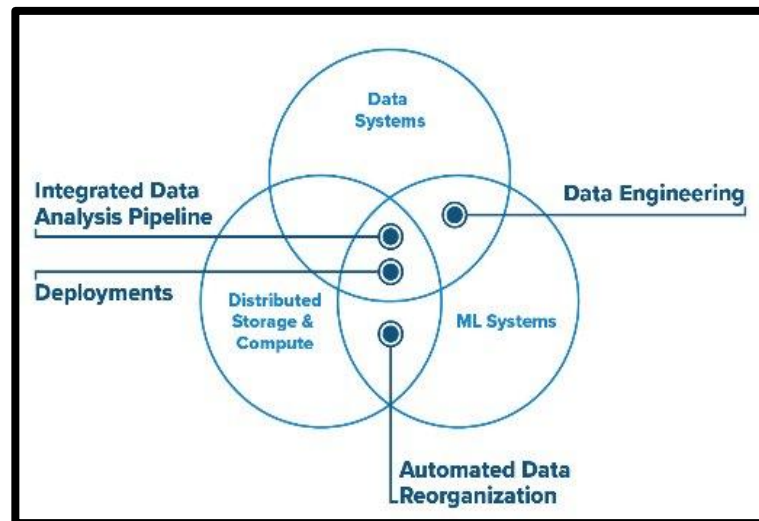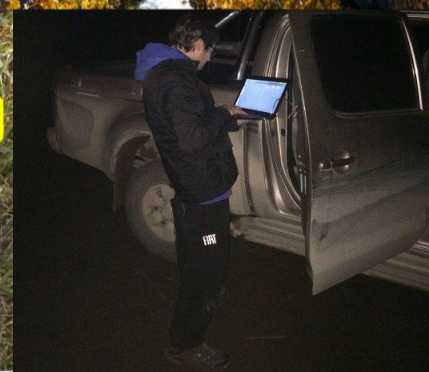| Area Manager | Researchers | Data Scientists |
|---|---|---|
|  |   |   |
| Lucas Iacono | Shafaq Siddiqi – Mark Dokter | Alexander Hiebl – Lorenz Dirry |

# Our research

———

- **Data Platforms and Cloud**
  - Adapt and tune data management and ML systems for domain-specific data platforms, federated learning, and hybrid cloud-edge.
- **Integrated Data Analysis Pipelines**
  - Simplify the full utilization of heterogeneous hardware, and adopt this infrastructure in related multi-firm projects
- **Automatic Data Reorganization**
  - Address the increasing redundancy in complex data science workflows for data preparation and cleaning, data augmentation, feature engineering, hyper-parameter optimization, and model training.

XBee Radio

AA Batteries

Sensirion
SHT15

Arduino PRO

POWERED BY

# Course Organization

# Basic Course Organization

———

- **Team**
  - Lecturers:
    - Dr. M.Sc. Shafaq Siddiqi, ISDS – Know Center Research
    - Lucas Iacono PhD., ISDS – Know Center Research
  - Tutor:
    - Saiful Islam, ISDS
- **Language**
  - Classes and slides: **English**
  - Communication and examination: **English**

# Basic Course Organization

———

- **Course Format**

  - **VU 2/1**, **5 ECTS** (2 x 1.5 ECTS + 1x2 ECTS), bachelor/master (+info next class)
  - **Weekly lectures** (Fri 3pm, HS i5 + Webex, including Q&A), attendance **optional**
  - **Mandatory exercises or programming project** (2 ECTS)
  - **Recommended papers** for additional reading on your own

- **Prerequisites**

  - **Preferred:** course Data Management / Databases (a very good boot system)
  - **Sufficient:** basic understanding of SQL / Data Management / Relational Algebra
  - **Basic programming skills** (**Python**, Java, R)

# Basic Course Organization

———

- Video Recording
  - Link in **TUbe & TeachCenter**
  - **Optional** attendance
  - **Hybrid,** in-person but video-recorded lectures
  - **HS i5** and **Webex:**
    - https://tugraz.webex.com/meet/shafaq.siddiqi (1st module)
    - https://tugraz.webex.com/meet/lucas.iacono (2nd module)

# and more...course organization

———

- Website
  - https://lucasiacono.github.io/ws2024_2025_dia.html
  - All course material (lecture slides) and dates
- Video Recording Lectures (TUbe)
- Communication
  - Informal language (Lucas and Shafaq)
  - Ask for feedback!!! (unclear content, missing background)
  - Newsgroup: N/A – email is fine, TeachCenter forum for discussions
  - Office hours: by appointment or after lecture
- Exam
  - **Completed exercises or project**
  - **Final written exam** (oral exam if <25 students take the exam and for Erasmus students)
  - **Grading** (30% project/exercises completion, 70% exam)

# and more...more...course organization

___

- **Course Applicability**
  - **Bachelor** program computer science (CS)
  - **Bachelor** program software engineering & management (SEM)
  - **Master** programs CS and SEM
    - Catalog Data Science: compulsory course in major/minor
  - **Free subject course** in any other study program or UNI

# Course Motivation and Goals

# Data Sources and Heterogeneity

— — —

**Some important concepts and keywords**

**Integration**

Process of combining and harmonizing data from multiple sources into a unified, coherent format that can be put to use for various analytical, operational and decision-making purposes.

IBM

Problem of combining data residing at different sources, and providing the user with a unified view of these data.

Lenzerini, M. (2002, June). Data integration: A theoretical perspective. Proceedings of the twenty-first ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems (pp. 233-246).

# Data Sources and Heterogeneity

— — —

**Some important concepts and keywords**

Homogeneity

Similarity

(e.g. formats, sources, and structures)

(e.g. HW architecture, OS, Communication Protocols)
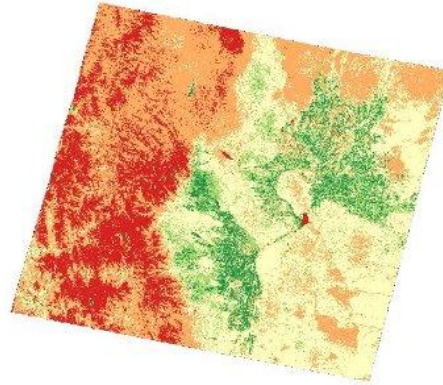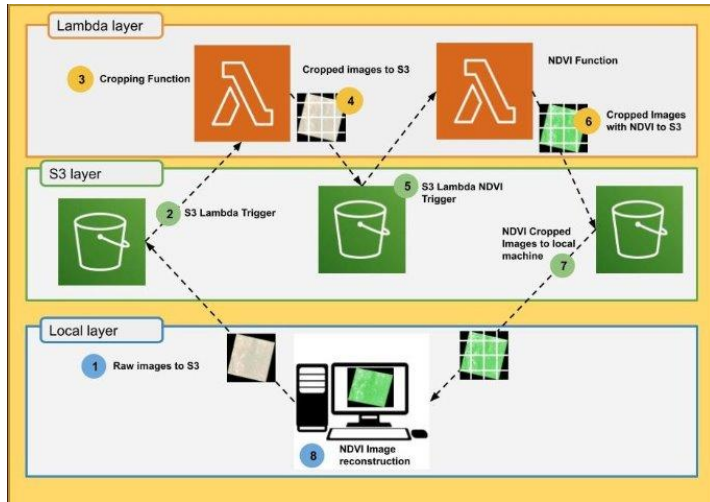
Heterogeneity

Dissimilarity

Data: CSV vs XLSX, %YY %DD %MM – MM/YY/DD , TS VS Social Media

HW: Local vs Cloud-based, ARM vs Intel

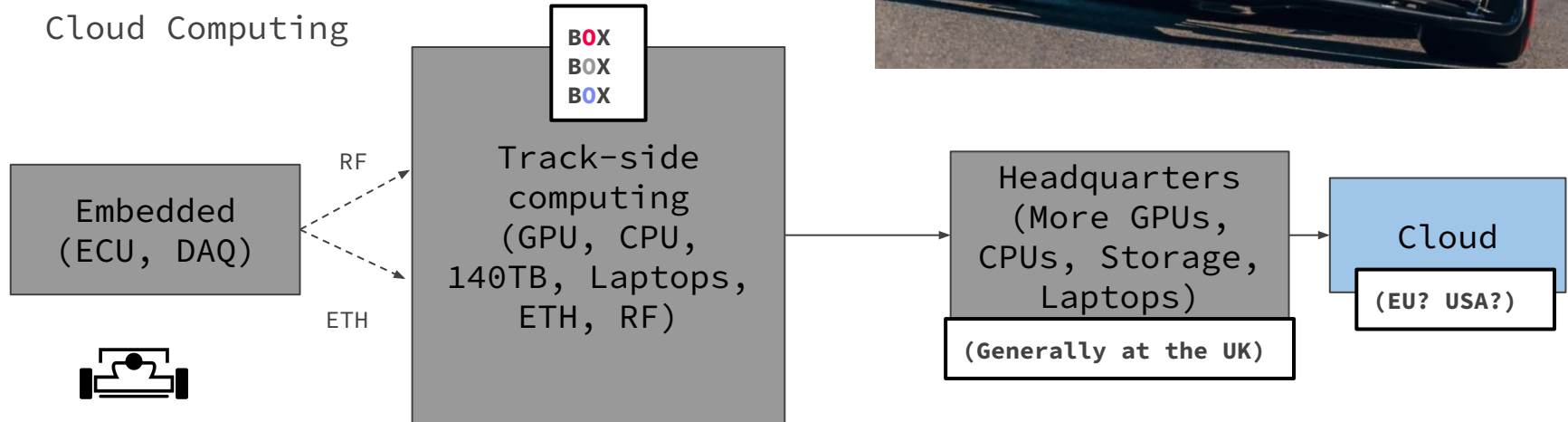# Data Sources and Heterogeneity: IT Infrastructure

Computing NDVI at the Cloud



Iacono, L., Pacios, D., & Vázquez-Poletti, J. L. (2023). SNDVI: A new scalable serverless framework to compute NDVI. Frontiers in High Performance Computing, 1, 1151530.

# Data Sources and Heterogeneity: Heterogeneous IT Infrastructure

---

- ECU and Car Systems
- Telemetry
- Laptops
- Track-side computing
- Headquarters computing
- Cloud Computing



```
Embedded          RF      Track-side
(ECU, DAQ)  ┄┄┄┄►  computing         ────►  Headquarters        ────►  Cloud
            ┄┄┄┄►  (GPU, CPU,               (More GPUs,                (EU? USA?)
              ETH   140TB, Laptops,          CPUs, Storage,
                    ETH, RF)                 Laptops)
                                             (Generally at the UK)
```

BOX
BOX
BOX

# Multi-modal data (Agriculture)

---

**Let's see if we can identify the heterogeneity**

- Structure
- Format
- Type
- Storage

**How we can integrate them?**
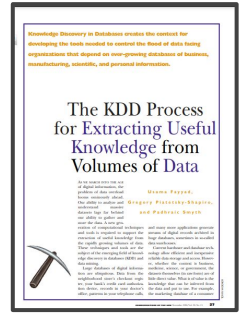
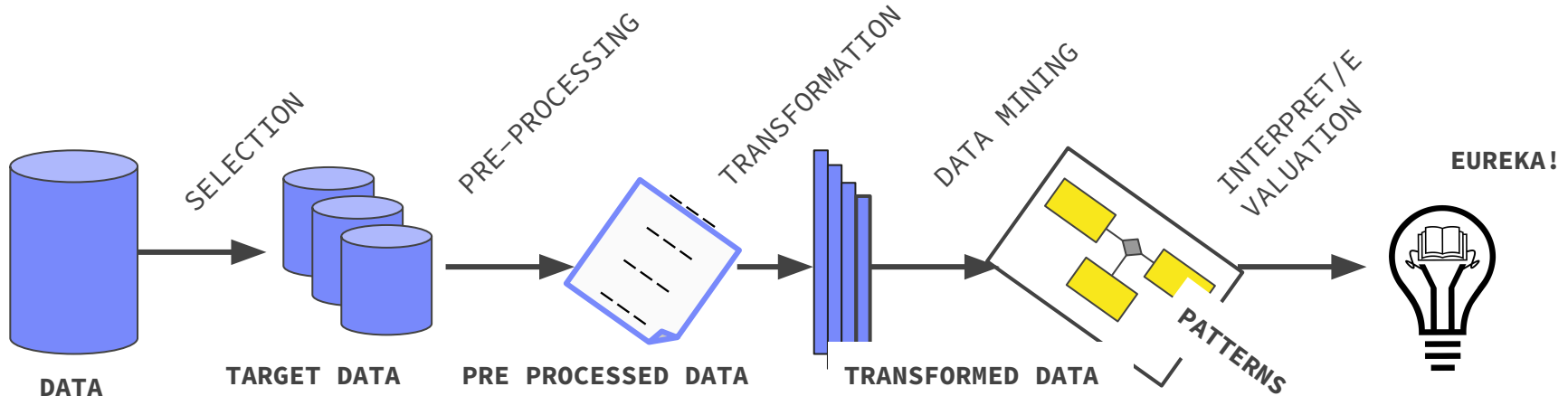# From data to models

———

# From KDD to the AI Lifecycle

— — —

**Classic KDD (Knowledge Discovery in Databases)**

**Descriptive** (association rules, clustering) and **predictive**

Fayyad, U., Piatetsky-Shapiro, G., & Smyth, P. (1996). The KDD process for extracting useful knowledge from volumes of data. *Communications of the ACM, 39*(11), 27-34.



DATA → SELECTION → TARGET DATA → PRE-PROCESSING → PRE PROCESSED DATA → TRANSFORMATION → TRANSFORMED DATA → DATA MINING → PATTERNS → INTERPRET/EVALUATION → EUREKA!
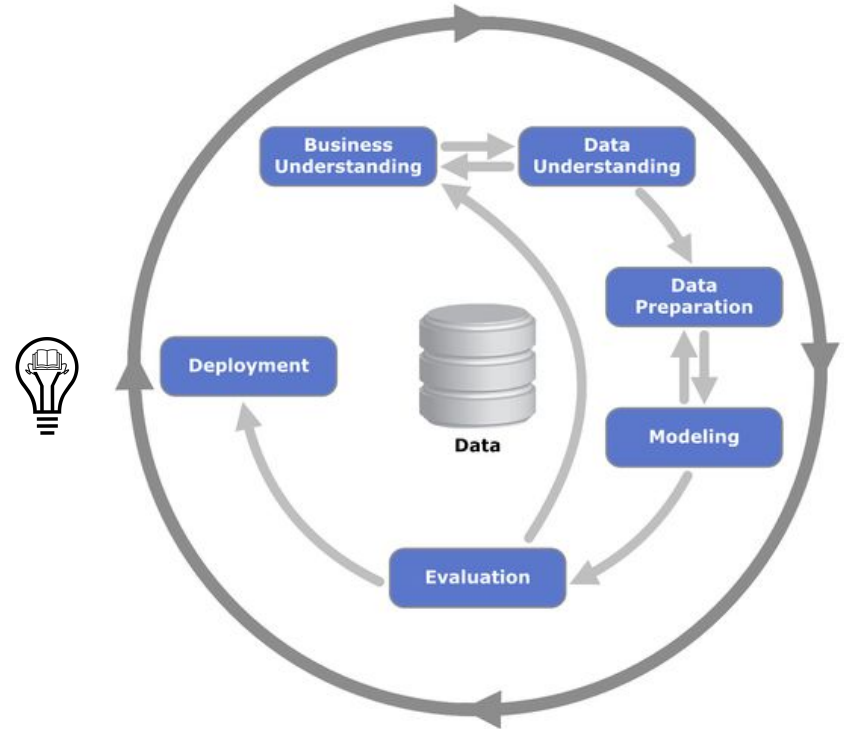
# From KDD to the AI Lifecycle

———

**CRISP-DM (Cross-Industry Standard Process for Data Mining)**

**What's new?** Business Understanding and Deployment **(A business perspective)**



Source: Statistik Dresden

# From KDD to the AI Lifecycle
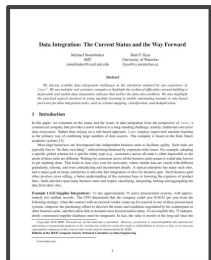
— — —

**AI Lifecycle**



Problem Definition → Data Acquisition and Preparation → Model Development & Training → Model Evaluation & Refinement → Model Deployment → MLOps: Model Monitoring & Maintenance → Problem Definition

# The 80% Argument

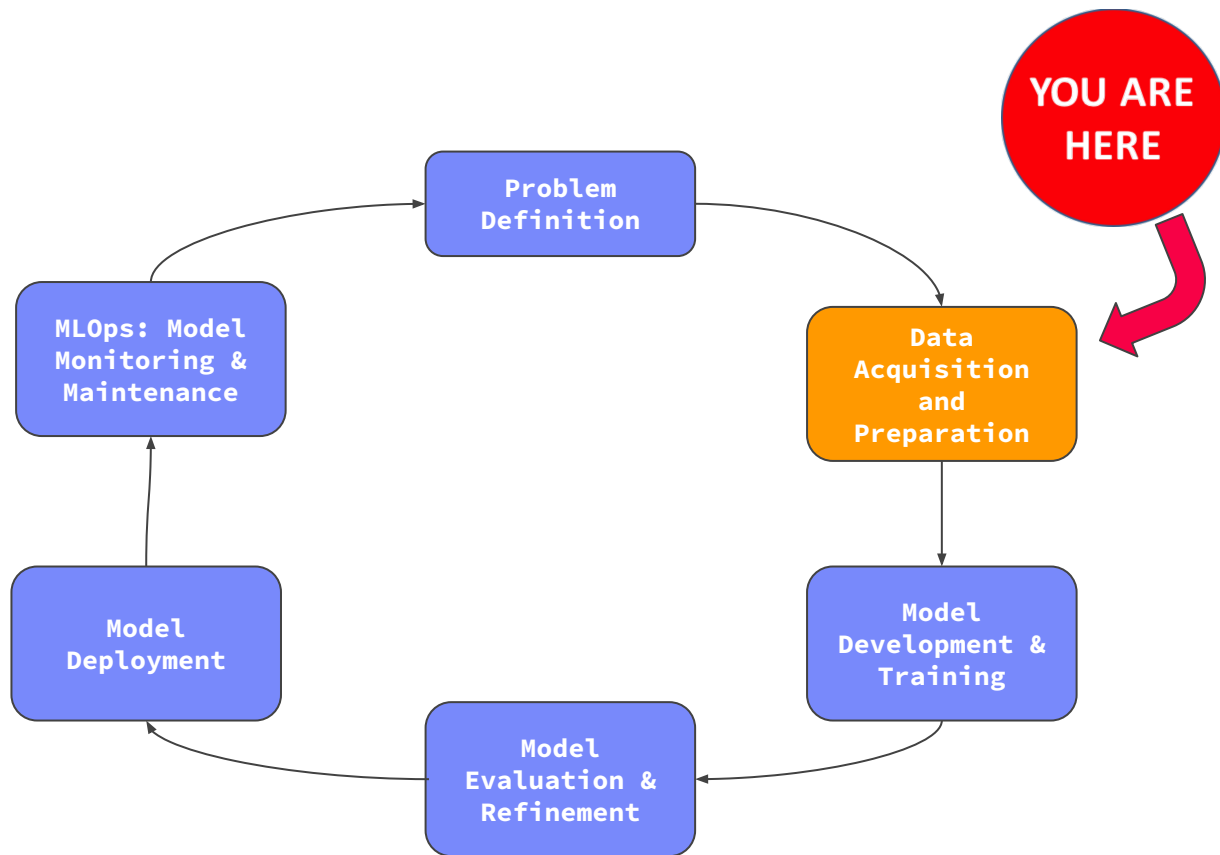Data scientists spend **80-90%** time on finding, integrating, cleaning datasets

Stonebraker, M., & Ilyas, I. F. (2018). Data Integration: The Current Status and the Way Forward. IEEE Data Eng. Bull., 41(2), 3-9.

YOU ARE HERE

Problem Definition

MLOps: Model Monitoring & Maintenance

Data Acquisition and Preparation

Model Deployment

Model Development & Training

Model Evaluation & Refinement

# Some open topics in the AI lifecycle

———

- **Multi-party data sharing.**
  - How to share data between organizations in multi-party and cross-country projects?
- **Dataspaces technologies.**
  - Data spaces have been conceptualized to guarantee sovereign data sharing in multi-party environments by the signature of digital agreements for the use and distribution of data.

# Some open topics in the AI lifecycle

———

- **Adaptability** Some of the most time-consuming processes in the AI lifecycle are detecting failures and decreases in data quality, such processes requires human intervention.
  - **Automated mechanisms** to predict and detect failures, avoiding the intervention of pipeline managers and flow guardians.
  - **Real-time data quality monitoring**. If there is a detection of a data inconsistency, the DQM module can request data resubmission, avoiding the involvement of humans.

# Some open topics in the AI lifecycle

———

- **Improvement of data traceability.** Some components of the AI lifecycle can experience errors. These errors can lead to data loss, and it is difficult to identify at which point in the pipeline the error occurred.
  - **Pipelines for AI need to automatically track the quality and reliability of the data and metadata obtained in each of their processes.**

# Course Goals

- Major data integration architectures
- Key techniques for data integration and cleaning
- Methods for large-scale data storage and analysis

# Learning Units and Projects

# Part A

Data Integration and Preparation

- LU1. Data Integration Architectures
  - Introduction and Overview [Oct 11]
  - Data Warehousing, ETL, and SQL/OLAP [Oct 18]
  - Message-oriented Middleware, EAI, and Replication [Oct 25]

— — —

# Part A

Data Integration and Preparation

- LU2. Key Integration Techniques
  - Schema Matching and Mapping [Nov 08]
  - Entity Linking and Deduplication [Nov 15]
  - Data Cleaning and Data Fusion [Nov 22]

---

# Part B

Large-Scale Data Management & Analysis

- LU3. Cloud Computing
    - Cloud Computing Fundamentals [Nov 29]
    - Cloud Resource Management and Scheduling [Dec 06]
    - Distributed Data Storage[Dec 13]

# Part B

Large-Scale Data Management & Analysis

- LU4. Large-Scale Data Analysis
  - Distributed, Data-Parallel Computation [Jan 10]
  - Distributed Stream Processing [Jan 17]
  - Distributed Machine Learning Systems [Jan 24]

# Overview Projects or Exercises

# Overview Projects or Exercises

— — —

- **Team**
  - **1 – 3 person teams (w/clearly separated responsibilities)**
- **Objectives**
  - Non-trivial programming project in DIA context **(2 ECTS  50 hours)**
  - **Exercise:** Data engineering and ML pipeline
    - Data cleaning and integration of multi-modal data sources
    - ML model training and evaluation
  - **Optional:** Open source contribution to Apache SystemDS or DAPHNE project **(from HW to high-level scripting)**
    - https://github.com/apache/systemds
    - https://github.com/daphne-eu/daphne
- **Timeline**
  - Oct 25: Exercise Description
  - Jan 17: Final project/exercise deadline

# Summary and Q&A

# Summary and Q&A

———

- **Course Goals**
  - Major data integration architectures
  - Key techniques for data integration and cleaning
  - Methods for large-scale data storage and analysis
- **Next Lectures**
  - Data Warehousing, ETL, and SQL/OLAP [Oct 18]
  - Message-oriented Middleware, EAI, and Replication [Oct 25]

# Vielen Dank!