# Introduction to Statistical Data Analysis and Machine Learning

Matteo Biagetti & Tommaso Rodani

Trieste, 22/10/2024

# Outline

- Why a data analysis course at MDMC?
- Objectives of the course
- Topics of the course
- Your data

# Why data analysis at MDMC?

# Why data analysis at MDMC?

Research Lab



Generated with AI

# Why data analysis at MDMC?



**FAIRification**: data and metadata acquisition procedure compliant with FAIR principles

## Research Lab



Generated with AI

# Why data analysis at MDMC?



**Automatisation**: Unsupervised data and metadata handling and pre-processing

Research Lab

**FAIRification**: data and metadata acquisition procedure compliant with FAIR principles

Generated with AI

# Why data analysis at MDMC?



**Automatisation**: Unsupervised data and metadata handling and pre-processing

**Data-driven approach:** Innovation through interdisciplinary research

WikiMedia Commons

Research Lab

Generated with AI

**FAIRification**: data and metadata acquisition procedure compliant with FAIR principles

# Why data analysis at MDMC?

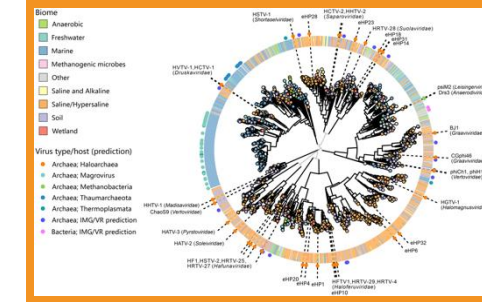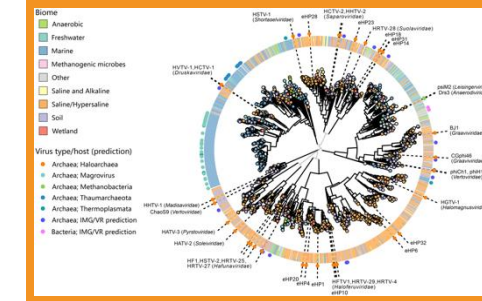**Automatisation**: Unsupervised data and metadata handling and pre-processing

**Data-driven approach:** Innovation through interdisciplinary research

FAIR Research Lab

WikiMedia Commons

**FAIRification**: data and metadata acquisition procedure compliant with FAIR principles

**Infrastructure-Augmentation:** Data and metadata analysis suggests improvements to the lab itself

Generated with AI

# Objectives

Suppose you have spent some effort to build a data workflow in your lab so that data and meta data are managed in a proper (FAIR) way.

## Now what?

- What type of information do these data and metadata contain?  **Problem Identification**

- Can they be preprocessed to simplify information extraction?  **Problem Investigation**

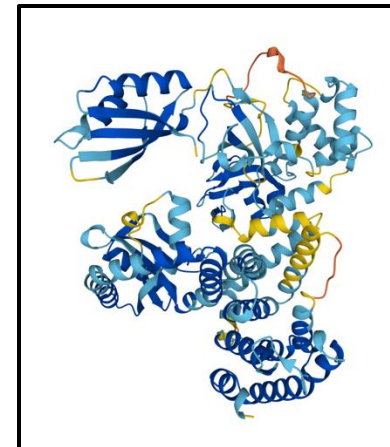- What tools can I use to extract information?  **Problem Solution**

# Topics

1. REGRESSION

Given an input-output pair, I want to learn the relation between the two

- To predict a new output given a new input
- To learn the physical, biological, chemical, etc, law between input and output
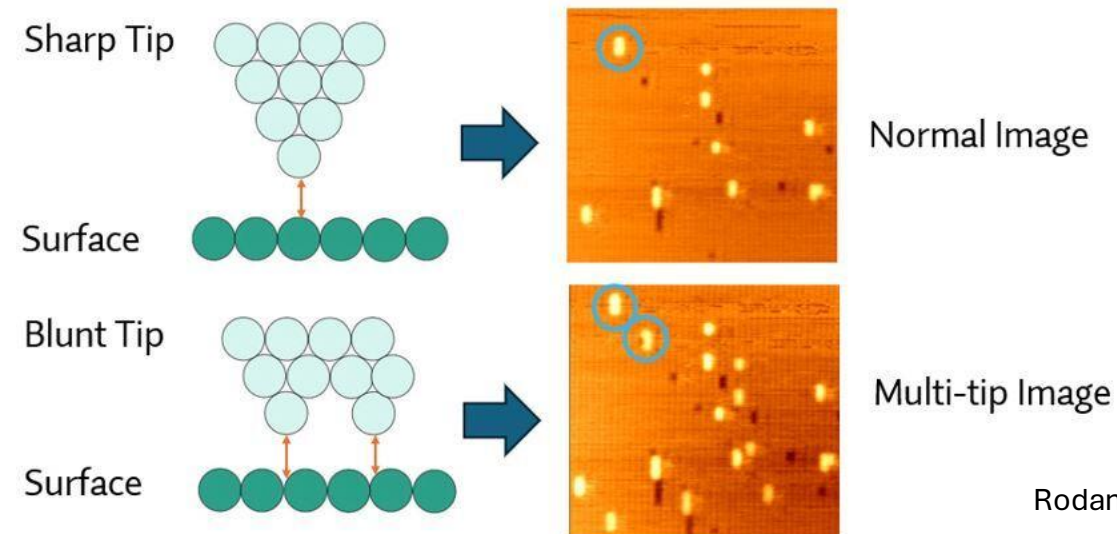- To generate a new pair
- ...



Braglia et al. 2017



AlphaFold (DeepMind)
WikiMedia Commons

# Topics

## 2. CLASSIFICATION

Given set of data with given properties, I want to classify them based on these properties

- To single out specific properties
- To filter data for purification
- …



Rodani et al 2024

# Topics

3. UNSUPERVISED LEARNING

Given set of data, I want to discover its properties

- To discover relations among quantities (new physical, biological, chemical laws)
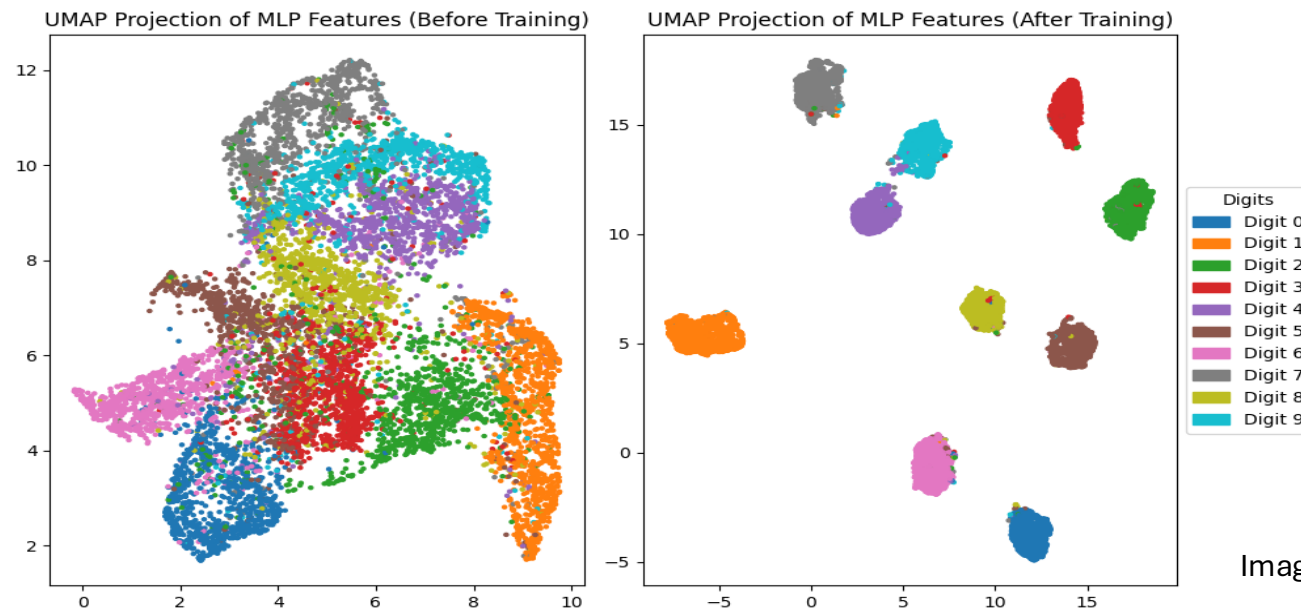- Reduce dimension of dataset
- ...



Image from T.Rodani

# Your Data!

nffa-di

Nano Foundries Fine Analysis
Digital Infrastructure

prp

Pathogen Readiness Platform
for CERIC-ERIC Upgrade