

The background of the slide is a dark blue field filled with a complex network of glowing blue lines and dots, resembling a data network or a molecular structure. The lines connect various points, some of which are highlighted with larger, brighter blue circles.

# Data Engineering SS25

## *Week 1 – Course Organization*

Dr. Sucheta Ghosh



HOCH  
SCHULE  
OFFEN  
BURG



- Introduction
- Course layout
- Course dates
- Exam
- Lab
- ...

- Dr. Sucheta Ghosh
  - Researcher at IWR University of Heidelberg
  - 
  - Lecturer in Data Engineering, Software Eng. Database Programming

## Planned: 12+ blocks

each with 2 SWS lecture + 2 SWS lab



- *Lecture*
  - Thursdays (12:00)
  - Room B122
- *Lab*
  - Thursdays
  - Room B106
  -



- Contact

- *[sucheta.ghosh@lehrbeauftrag.hs-offenburg.de](mailto:sucheta.ghosh@lehrbeauftrag.hs-offenburg.de)*
- For technical questions: please use the anonymous Moodle forums!



Please sign up to the course Moodle:

<https://elearning.hs-offenburg.de/moodle/course/view.php?id=6852>

There is a technical problem, I have contacted the HelpDesk

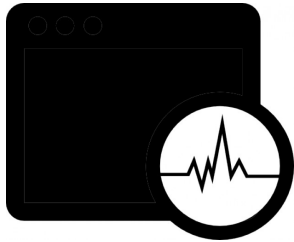


Code + Lab materials on GitHub:

<https://github.com/ghoshsucheta/DataEng25>

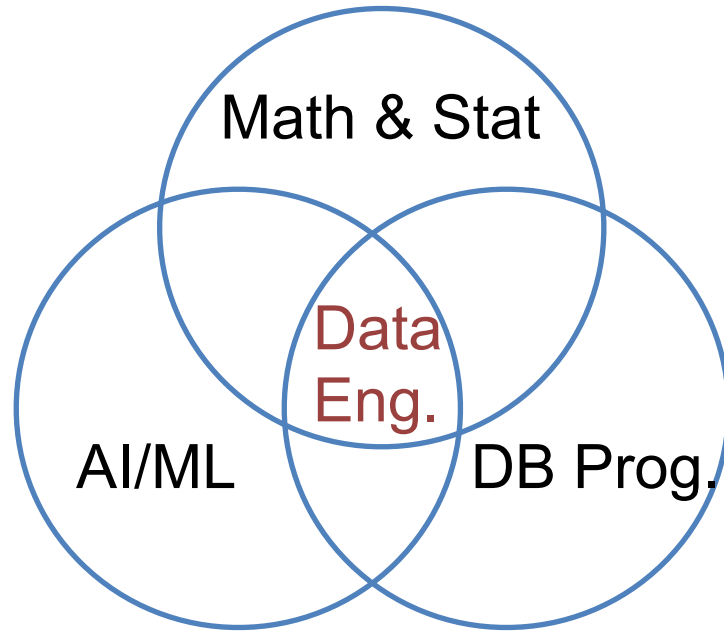


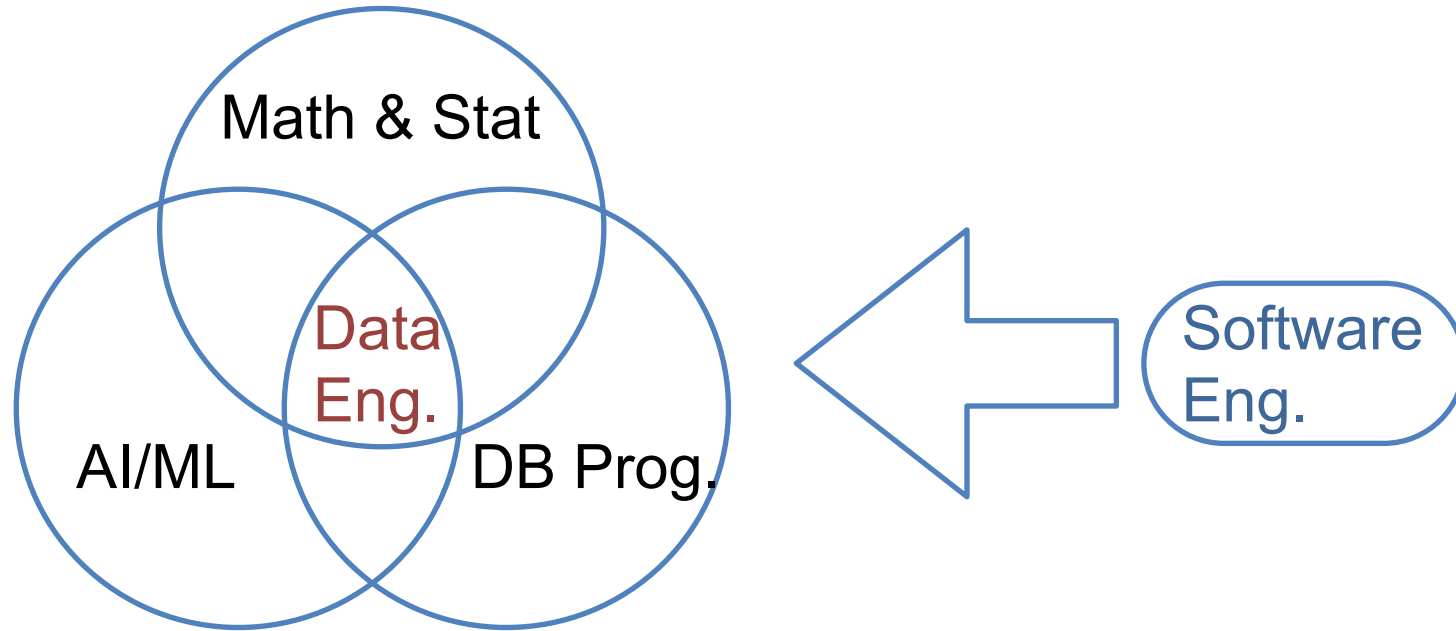
- 60/90min written exam (details later)



- Lab exercises + home work
  - Mandatory !
  - new exercises every Thursday
  - submit via Moodle by next Wednesday 11.59pm

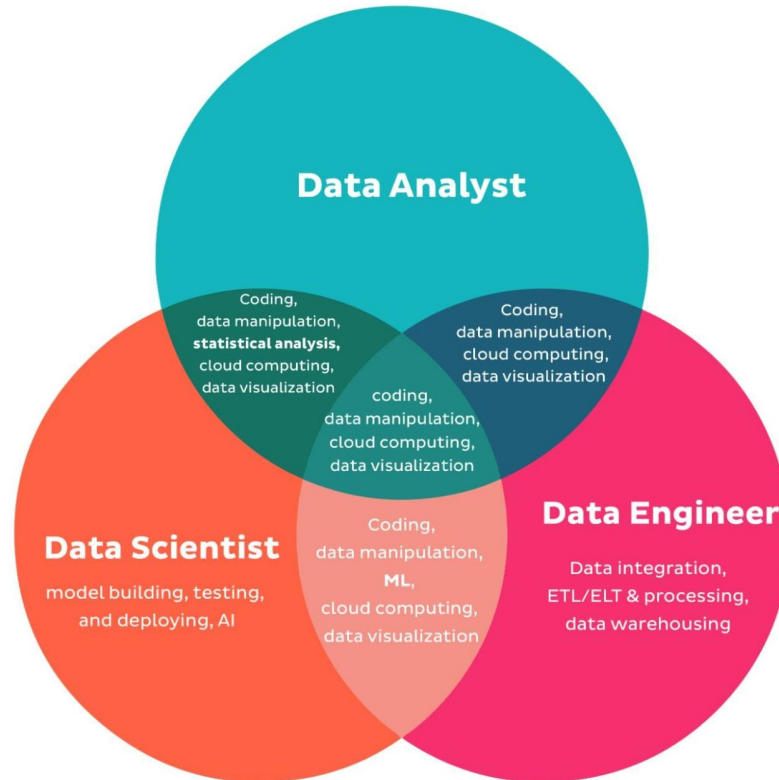






# Data Engineering vs. Data Analysis vs. Data Scientist

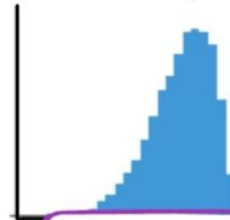
## REQUIRED SKILLS



# Data Engineering vs. Data Analysis vs. Data Scientist

| Dimension        | Data Analyst   | Data Scientist   | Data Engineer  |
|------------------|--|--|--|
| Focus            | Data analysis  | Predictive models  | Data infrastructure  |
| Skills           | Statistical analysis, programming in R and Python              | Machine learning, programming in R and Python                                      | Database technologies, programming languages                                     |
| Responsibilities | Collecting, processing, analyzing data, making recommendations | Designing and developing predictive models, providing insights and recommendations | Designing, building, and maintaining data infrastructure, ensuring data accuracy |
| Technologies     | SQL, R, Python   | Python, R  | SQL, NoSQL, Hadoop, Spark  |

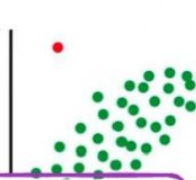
## EDA Exploratory Data Analysis



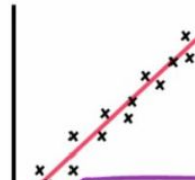
Data distribution

|   | F        | G        | H        | I        | J        |
|---|----------|----------|----------|----------|----------|
| A | 0.620576 | 0.140053 | 1.352728 | NaN      | 0.808078 |
| B | NaN      | 0.526829 | NaN      | NaN      | 0.170902 |
| C | NaN      | 0.458827 | 1.406713 | 0.071119 | NaN      |
| D | NaN      | 2.307197 | NaN      | NaN      | NaN      |
| E | 0.203402 | 0.259913 | NaN      | 0.505811 | 1.516755 |

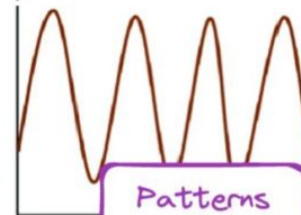
Missing data



Outliers



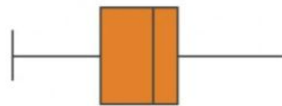
Correlation



Patterns

```
Cust_No          int64
Cust_Name        object
Product_id       int64
Product_cost     float64
Purchase_Date    datetime64[ns]
dtype: object
```

Data types















Data visualization



Data quality

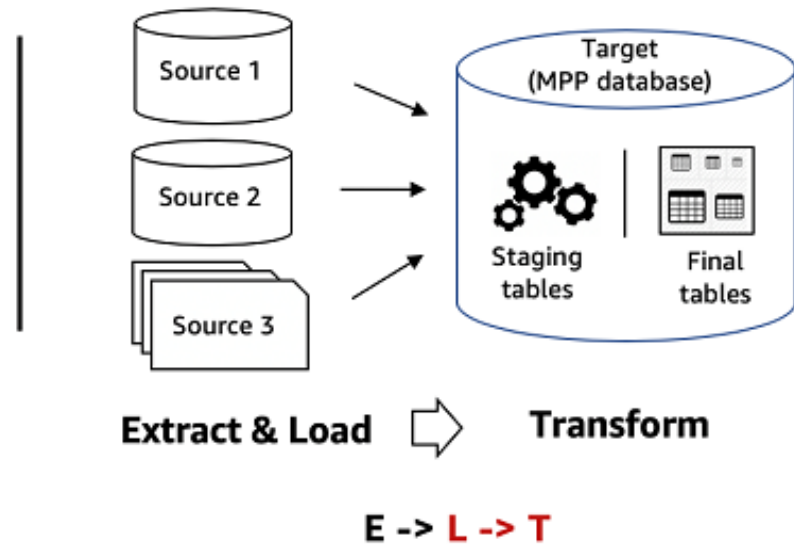
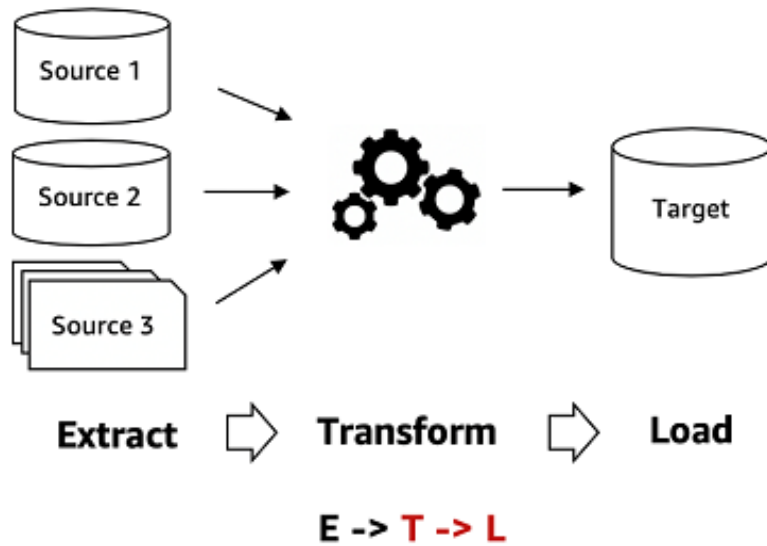
## Data Visualization: Key differences between approaches

|                            | Exploratory  | Explanatory   |
|----------------------------|--|---|
| <b>Goal</b>                | Understand             | Communicate                        |
| <b>Audience</b>            | You                    | Other people                       |
| <b>Data Familiarity</b>    | Very Familiar (You)    | Less familiar (Others)             |
| <b>Visualization Focus</b> | Flexibility and speed  | Simplicity, clarity, and cohesion  |
| <b>Narrative</b>           | Unknown                | Known                              |
| <b>Outcome</b>             | Insight               | Action                            |

Effectivedatastorytelling.com

Next Part of this presentation is for practical part,  
now the rest of the lecture today will be borrowed  
from Prof. Dr.-Ing. Janis Keuper (Last Sem. week 1)

# ETL vs ELT Data Engineering

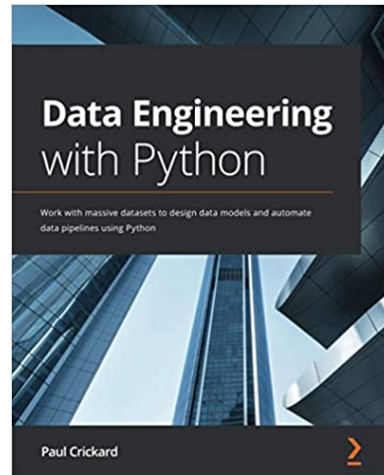
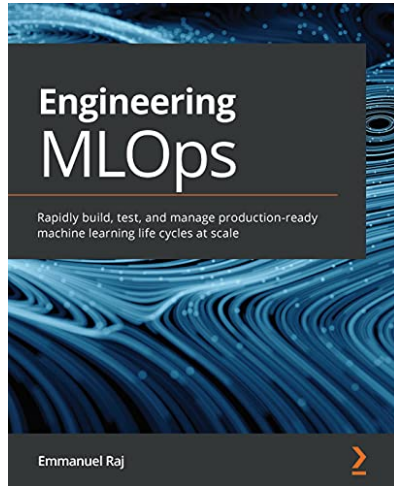




A Tutorial:

<https://www.neonscience.org/resources/learning-hub/tutorials/about-hdf5>

For now I suggest



**Non of the books cover this course completely, it is NOT mandatory to have any of them.**

[1] free icons taken from <https://www.flaticon.com>