



Software Practicals

Winter Semester 2023/24

Data Science Group
Heidelberg University
October 18, 2023

Slides Online



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386



The slides are available on our webpage
<https://ds.ifi.uni-heidelberg.de/teaching/current/>



Organization

Outline



- Overview of topics (today)
 - Send application for a topic until **Monday, October 23, 1pm**
 - Assignment of topics by October 25
- First milestone (before Christmas break)
 - Prototype / part of software
 - Summary of research (literature and related systems/tools)
 - Further milestones in agreement with supervisor
- End of practical (mid/end February)
 - Code has to be in local Gitlab of the database group
 - Presentation / demo of practical and software (10-12 minutes)
 - Report / documentation as Gitlab document (README.md)

Application



- Apply directly to supervisor via mail
 - Program of study, semester of study, matriculation number
 - List relevant course experience, including course grades
 - List other experience:
 - Side projects you are working on
 - “Anwendungsgebiet” / Application Field
 - Job and project experience
 - Send your tentative schedule and milestones for the practical
 - Group work is not possible!
- It is recommended to apply for multiple topics (e.g., “top-3 list”)

Application is binding!

Don't apply if you don't want to do the practical!

Deadlines



- In general: biweekly meetings with supervisor
- Presentation: end of February 2024
- Report & Gitlab upload: end of February 2024
- No extension possible

Not finished = failed (grade 5,0)!

- Credit points (Leistungspunkte)
 - Beginners Practical (IAP, 2 CP + 4 FÜK) [Bachelor students]
 - workload: 180 h (~1 ½ days/week)
 - Advanced Practical (IFP, 8 CP)
 - workload: 240 h (~2 days/week)
- Grading based on
 - code (readability, structure, functionality; code in local Gitlab)
 - documentation (README.md, code comments, documentation in Gitlab)
 - commitment and self-reliance
 - cool ideas!!
- IMPORTANT
 - talk to / communicate with your advisor (at least biweekly meetings)

Supervisors



- Michael Gertz (MG)
gertz@informatik.uni-heidelberg.de
- Satya Almasian (SA)
almasian@informatik.uni-heidelberg.de
- Jayson Salazar (JS)
salazar@informatik.uni-heidelberg.de
- John Ziegler (JZ)
ziegler@informatik.uni-heidelberg.de
- Ashish Chouhan (AC)
chouhan@informatik.uni-heidelberg.de
- Nicolas Reuter (NR)
reuter@informatik.uni-heidelberg.de

Project Topics

AP = Advanced Topic

BP = Beginners Topic (for BSc students)

Overview of Topics

1. Quantity and Concept Extraction with ChatGPT, **BP** (Gertz/Almasian)
2. Form Extraction from OCR'd Documents, **AP** (Gertz)
3. Graph Library Benchmark, **AP** (Ziegler)
4. Visual Benchmark of Dimensionality Reduction for Big, Sparse Graphs **AP** (Salazar)
5. Acquisition, Analysis and NER on OPS Codes **AP/BP** (Salazar)
6. Package Integration for NER in Patient Records **BP/AP** (Salazar)
7. Unlocking Legal Insights, **BP** (Gertz/Chouhan)
8. Concept Exploration UI, **AP** (Chouhan)
9. Temporal Evolution of Legal Documents, **BP/AP** (Chouhan)
10. Table Structure Recognition with Ruling Lines, **AP/BP** (Reuter)
11. Creating a PDF Table Annotation Tool, **AP** (Reuter)
12. Creating a Domain-Aware PDF Table Corpus, **AP** (Reuter)

BP: Quantity and Concept Extraction with ChatGPT (MG/SA)



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

Given:

- Sentences from news articles tagged with quantity information

Tasks:

- Use ChatGPT to extract quantity information for a sentence, containing <value, unit, change, concept>
- “Apple hires 200 people”
→ value=200, unit=people, change=equal, concept=Apple.



Subtasks:

- Examine ChatGPT’s ability to standardize values, normalize unit, and find relevant nouns (concepts)
- Build a pipeline that gets sentences and outputs quantity information

Languages / Tools:

- Python

AP: Form Extraction from OCR'd Documents (MG)



Given:

- Collection OCR'd forms, e.g., receipts, bills, table like structure
- Popular libraries and tools to perform AI-based layout analysis

Tasks:

- Evaluate quality of different tools, frameworks and libraries
- Build framework that allows (visual) comparison of different tools

Subtasks:

- Familiarize yourself with [document layout analysis](#), [LayoutLM](#)

Languages / Tools:

- Python, Hugging Face

AP: Graph Library Benchmark (JZ)

Given:

- Various open-source graph libraries
- Past benchmarks exist, e.g., [here](#)

Task:

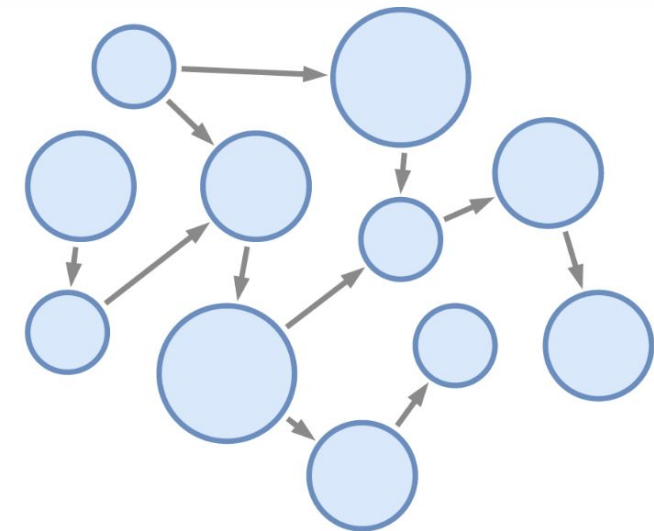
- Conduct performance benchmark of different graph libraries

Subtasks:

- Prepare benchmark dataset(s)
- Create “lab” environment
- Special focus on *network dynamics*

Languages / Tools:

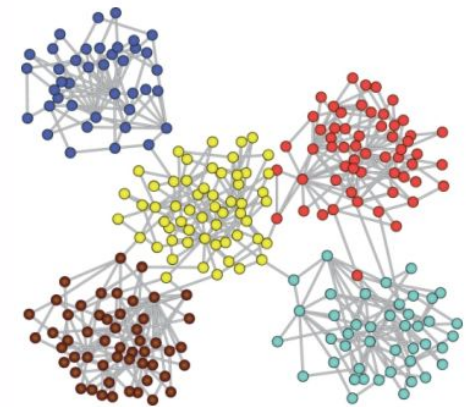
- Python, graph-tool, igraph, NetworkX, ...



AP: Visual DimRed Benchmark for Graphs AP (JS)

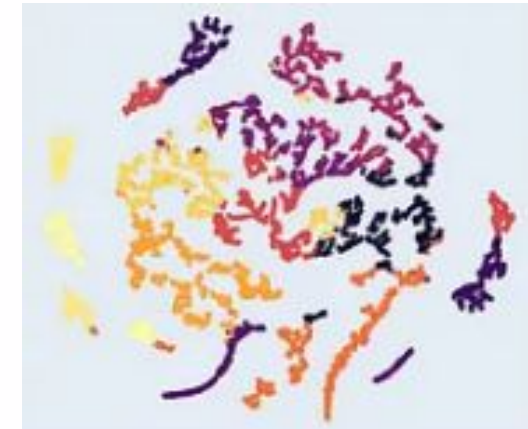
Given:

- Graphs can be built off almost any sort of text domain data, but they also grow **quickly** and are **sparse**.
- Representing, grouping and visualizing their labels as well as properties is of crucial importance in our group



Tasks:

- Generate a synthetic graph dataset based on given, well-defined properties from real-world data examples.
- Build an application that allows a user to poll graph data, brush(filter) it and visualize it based on two chosen algorithms (e.g. PCA, t-SNE and UMAP).



Languages / Tools:

- Python, Javascript (Cytoscape.js+Svelte), Apache AGE (or Neo4J)

BP(AP): Acquisition, Analysis and NER on OPS Codes (JS)



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

Given:

- Python library for Medical Thesaurus Correlation (MedKEET)
- Access to OPS-Database and related medical thesauri

Tasks:

- Acquire and analyze the structure and contents of the OPS dataset.
- Extend MedKEET to process and annotate (naively) OPS entries

Languages / Tools:

- Python, Pandas, Neo4J

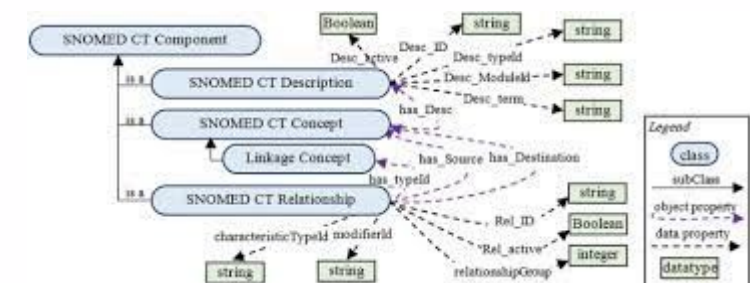


Federal Institute
for Drugs
and Medical Devices

Back in 2000 , **People Magazine** **PUBLISHER**
the time was a little more fashion-conscious , e

Now-a-days the prince mainly wears **navy** **COL**
double-breasted **DESIGN**) , **light blue** **COL**
pointed **DESIGN** **collars** **PART** , and **burg**

But who knows what the future holds ...
Duchess Kate **PERSON** did wear an **Alexan**
wedding **OCCASION** in the **fall of 2017** **SEA**





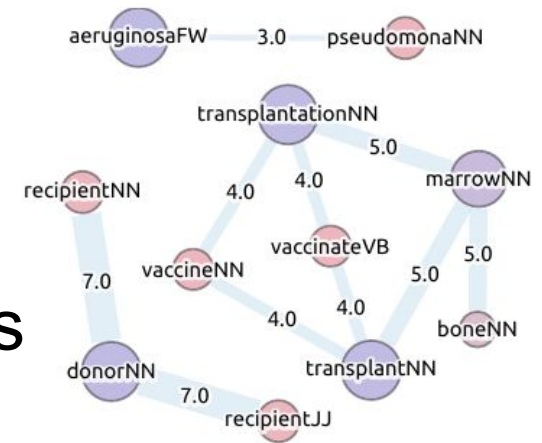
BP(AP): Package Integration for NER in Patient Records (JS)

Given:

- PyCoNet, a prototype package aimed towards the extraction of concept relationships from raw (medical) text.
- [FHIRPACK](#), an open source [FHIR](#) data Python processing toolkit
- Both developed jointly in the DBS and the UK-Essen

Tasks:

- Integrate and extend both packages so is-relationships (NER) present in raw, semi-structured patient records can be extracted and stored seamlessly



Languages / Tools:

- Python, PostgreSQL

BP: Unlocking Legal Insights (MG/AC)



Given:

- People post information about legal domain on LinkedIn (see, e.g., [Martin Ebers](#))

Tasks:

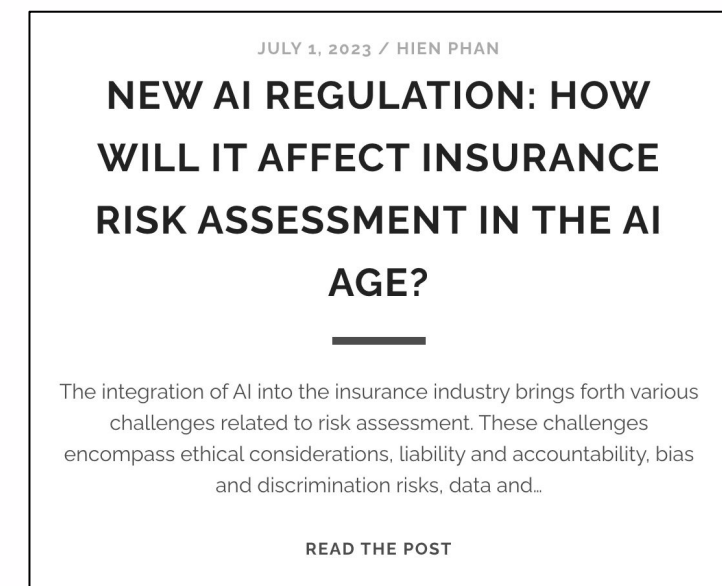
- Extract articles, documents, and posts related to legal domain and store it in [OpenSearch](#)

Subtasks:

- Realize backend to store extracted information using [OpenSearch](#)
- Analyze dataset and compute basic statistics

Languages / Tools:

- Python, [OpenSearch](#)



[Link](#)

AP: Concept Exploration UI (AC)



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

Given:

- API provides weighted word co-occurrence network for concept exploration in Pubmed abstracts

Tasks:

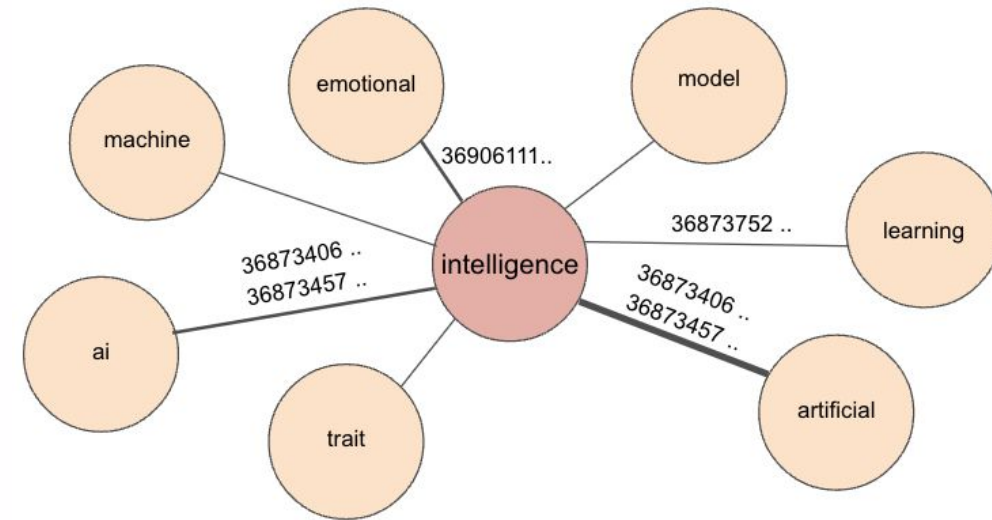
- Implementation of UI to explore concept

Subtasks:

- Handle API access
- Visualization of networks

Languages / Tools:

- Python, [Sveltekit](#), [eCharts](#), [Cytoscape.js](#)



Given:

- [Consolidated text](#) from EUR-Lex website (HTML)

Tasks:

- Fetch consolidated text information
- **AP:** develop frontend to search and explore evolution of regulations

► B ↓				AGREEMENT
				on the withdrawal of the United Kingdom of Great Britain and Northern Ireland from the European Union and the European Atomic Energy Community
				(OJ L 029 31.1.2020, p. 7)
Amended by:				
				Official Journal
				No page date
► M1 ↓				DECISION No 1/2020 OF THE JOINT COMMITTEE ESTABLISHED BY THE AGREEMENT ON THE WITHDRAWAL OF THE UNITED KINGDOM OF GREAT BRITAIN AND NORTHERN IRELAND FROM THE EUROPEAN UNION AND THE EUROPEAN ATOMIC ENERGY COMMUNITY of 12 June 2020
				L 225 53 14.7.2020
► M2 ↓				DECISION No 3/2020 OF THE JOINT COMMITTEE ESTABLISHED BY THE AGREEMENT ON THE WITHDRAWAL OF THE UNITED KINGDOM OF GREAT BRITAIN AND NORTHERN IRELAND FROM THE EUROPEAN UNION AND THE EUROPEAN ATOMIC ENERGY COMMUNITY of 17 December 2020
				L 443 3 30.12.2020



EUR-Lex
Access to European Union law

Subtasks:

- Design and implement document store based on [OpenSearch](#) and [PostgreSQL](#)

Languages / Tools:

- Python; [SPARQL](#); [OpenSearch](#); [PostgreSQL](#)

BP/AP: Table Structure Recognition with Ruling Lines (NR)



Given:

- [FinTabNet](#) dataset of annotated PDFs with Financial Report Tables

	2002		2001	
	High	Low	High	Low
1st Quarter	\$ 59.46	\$ 48.32	\$ 64.00	\$ 45.75
2nd Quarter	59.65	49.54	59.55	50.26
3rd Quarter	55.34	42.30	60.04	50.23
4th Quarter	44.37	34.14	53.61	42.30

Task:

- Implement a method that uses ruling lines to extract the structure of complex tables (mapping cells to the correct headers)

U.S. Pension Benefits			Non-U.S. Pension Benefits		
Fiscal Year			Fiscal Year		
2015	2014	2013	2015	2014	2013
\$ 104	\$ 107	\$ 104	\$ 60	\$ 54	\$ 43
105	97	94	33	29	27
(160)	(141)	(128)	(41)	(35)	(33)
—	1	(1)	—	1	1
65	85	71	12	11	8
\$ 114	\$ 149	\$ 140	\$ 64	\$ 60	\$ 46

Subtasks:

- Identify tables with rulings in the dataset and extract the ruling lines
- Recognize if a ruling belongs to a heading and assign it to it
- Use the position and size of the lines to extract the table structure

Languages / Tools:

- Python, [pdfplumber](#), (Docker)

AP: Creating a PDF Table Annotation Tool (NR)

Given:

- Reference annotations: [FinTabNet](#) dataset
- Possible starting point: table extraction tool from previous practical

Task:

- Create a tool with a simple UI for annotating (complex) tables

Subtasks:

- Extract the tables from PDFs ([camelot](#) and [tabula-py](#) can be used)
- Create a UI with which the extracted tables can be edited to get a proper annotation
- Convert the annotations to the target format (see FinTabNet)

Languages / Tools:

- Python, JavaScript (e.g. node.js, Svelte), Docker

A summary of future minimum lease payments under noncancelable operating leases with an initial or remaining term in excess of one year at May 31, 2017 is as follows (in millions):

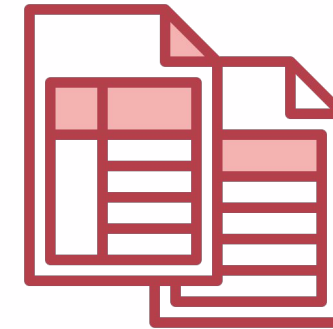
	Operating Leases	Aircraft and Related Equipment	Facilities and Other	Total Operating Leases
2018	2,441	398	2,047	4,886
2019	4,260	645	1,687	6,592
2020	4,381	261	1,670	6,312
2021	4,709	208	1,508	6,425
2022	4,540	184	1,351	6,075
Thereafter	3,019	175	1,844	5,038
Total	\$17,874	\$1,501	\$10,508	\$29,883

AP: Creating a Domain-Aware PDF Table Corpus (NR)



Given:

- A very large dataset of PDFs, e.g.
[CC-MAIN-2021-31-PDF-UNTRUNCATED](#)



Tasks:

- Create a subset of the dataset containing only PDFs with tables
- Classify the PDFs by domain/content, language (and table type)

Subtasks:

- Implement a pipeline for table detection, metadata extraction and PDF classification
- Basic data cleanup: find duplicates and spam/useless documents using the extracted data
- Optional: Implement a method to crawl new PDFs with tables

Languages / Tools:

- Python, PostgreSQL, (Docker, OpenSearch)

Slides Online



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386



The slides are available on our webpage
<https://ds.ifi.uni-heidelberg.de/teaching/current/>