



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

Harvesting German clinical Knowledge Graphs from pre-trained LMs

Natalia Minakova

MA Thesis Presentation

Prof. Dr. Anette Frank

Prof. Dr. ret. nat. Christoph Dieterich

Overview

- Knowledge Graphs
 - Motivation
 - Clinical Domain
- Prior Work
- Research questions
- Resources
- Experiments
- Results

Knowledge Graphs - Motivation

Definition:

“a graph of data intended to accumulate and convey knowledge of the real world, whose nodes represent entities of interest and whose edges represent relations between these entities” (Hogan et al. 2021)

- Effective way of representing information in a structured manner: a network of interconnected entities, attributes, and relationships

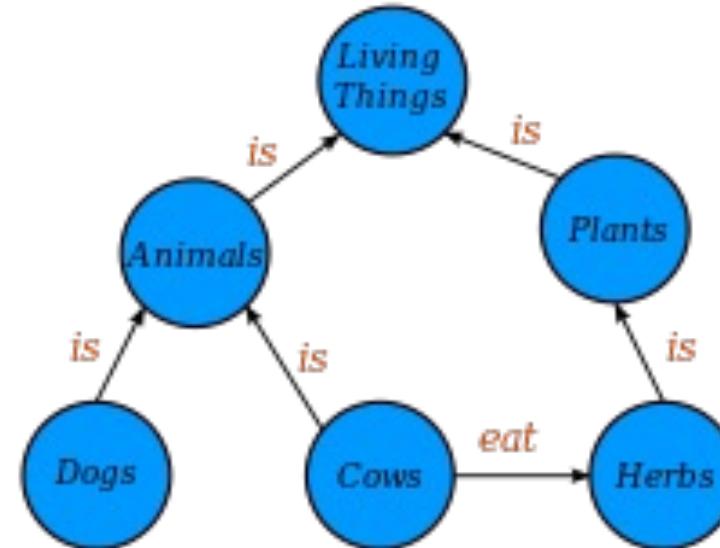


Figure 1. KG sample [1]

Knowledge Graphs in the Health Domain

Electronic Health Records contain massive quantities of unstructured patient data

Advantage of using KG in healthcare:

- Providing a unified view of patient's information
- Capturing the semantics and relationships between medical concepts leading to improvements in tasks like Natural Language Understanding
- Supporting clinical decisions:
 - fast, easy, and efficient access to health-related information, enables the physicians to make more informed decisions at the point of care
- Facilitating Precision Medicine and Personalized Healthcare
- Efficient Information Retrieval
 - more precise and context-aware search queries, for example, to access clinical guidelines

How can we construct a Knowledge Graph?

- Crowd sourcing
- Text mining pipelines to extract knowledge from text:
 - Entity extraction, coreference resolution, entity linking and relation extraction

What if there was another way?

Language Models as Knowledge Bases?

(Petroni et al. 2019)

- Can a LM be queried for relational data the same as a knowledge graph?
- Language models:
 - No schema engineering
 - No need for human annotation
 - Support a open set of queries

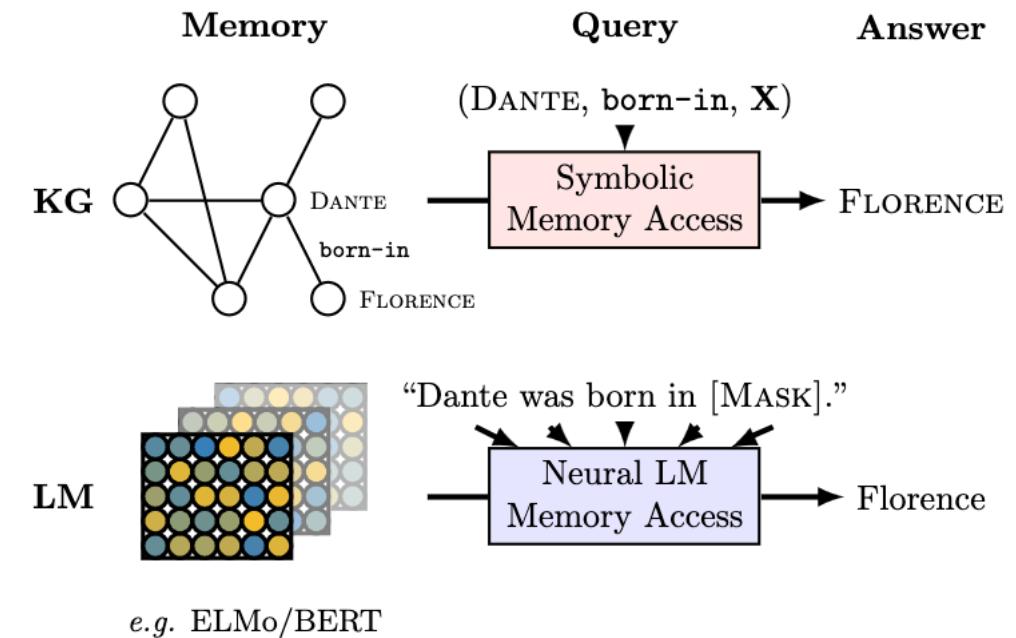


Fig 2. Querying KB and LM for factual knowledge^[1]

Crawling The Internal Knowledge-Base of Language Models (Cohen et al. 2023)

- Crawling procedure: starts from the seed entity and recursively expands it to expose additional facts

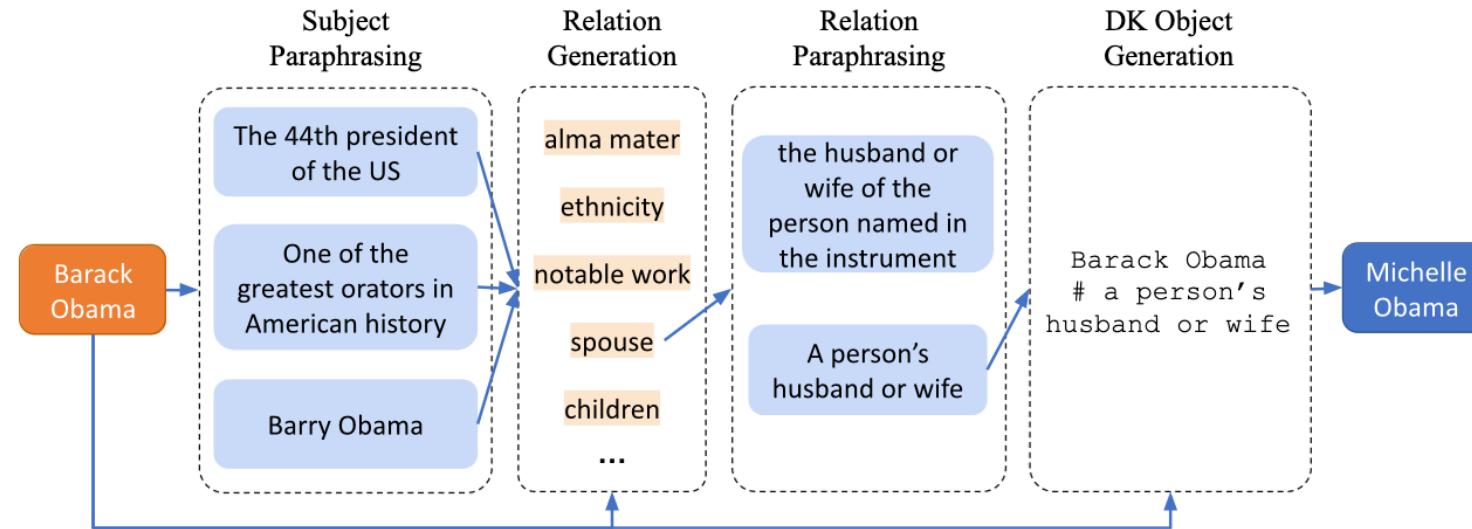


Fig 3. An illustration of the full method for crawling a subgraph [1]

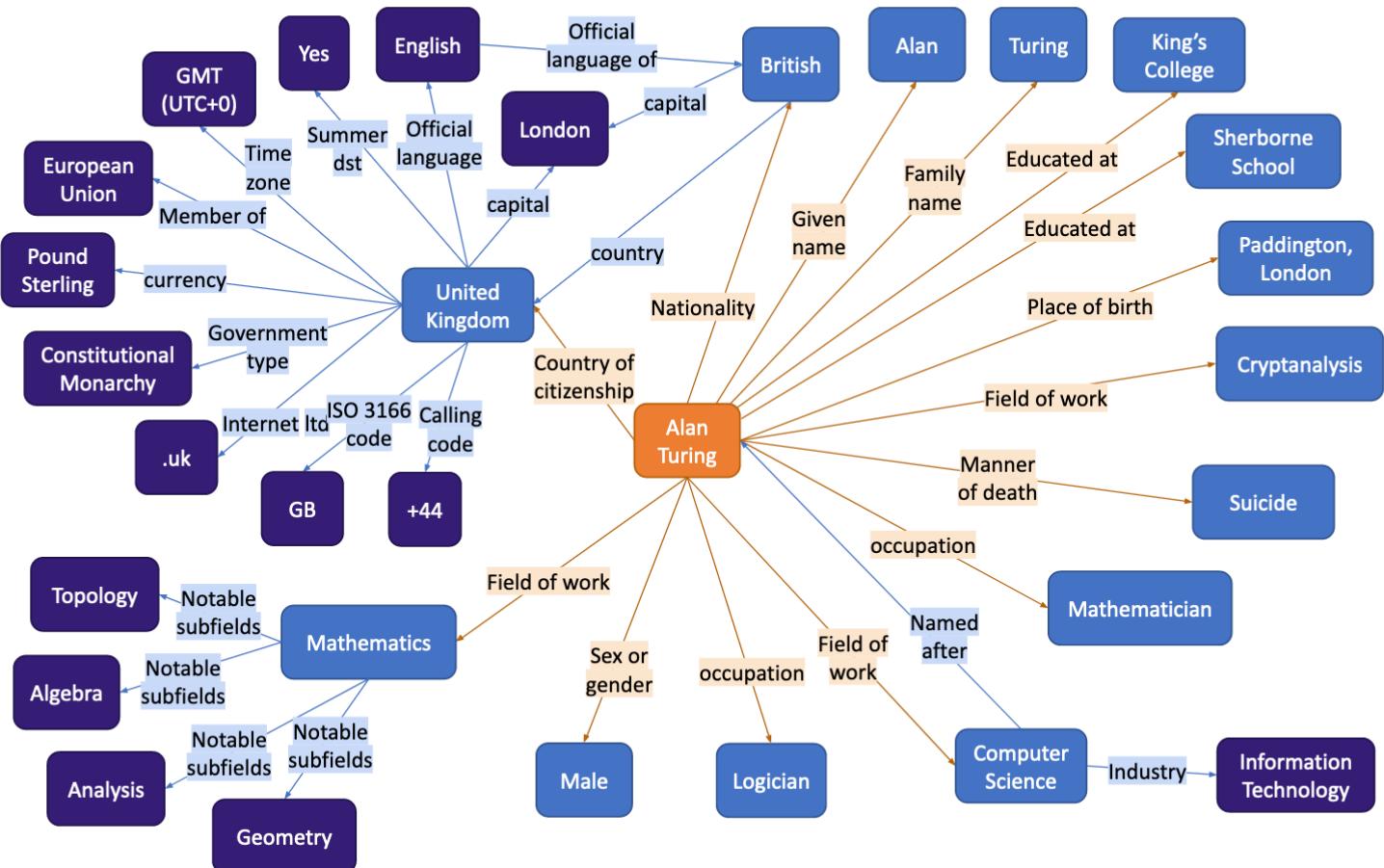


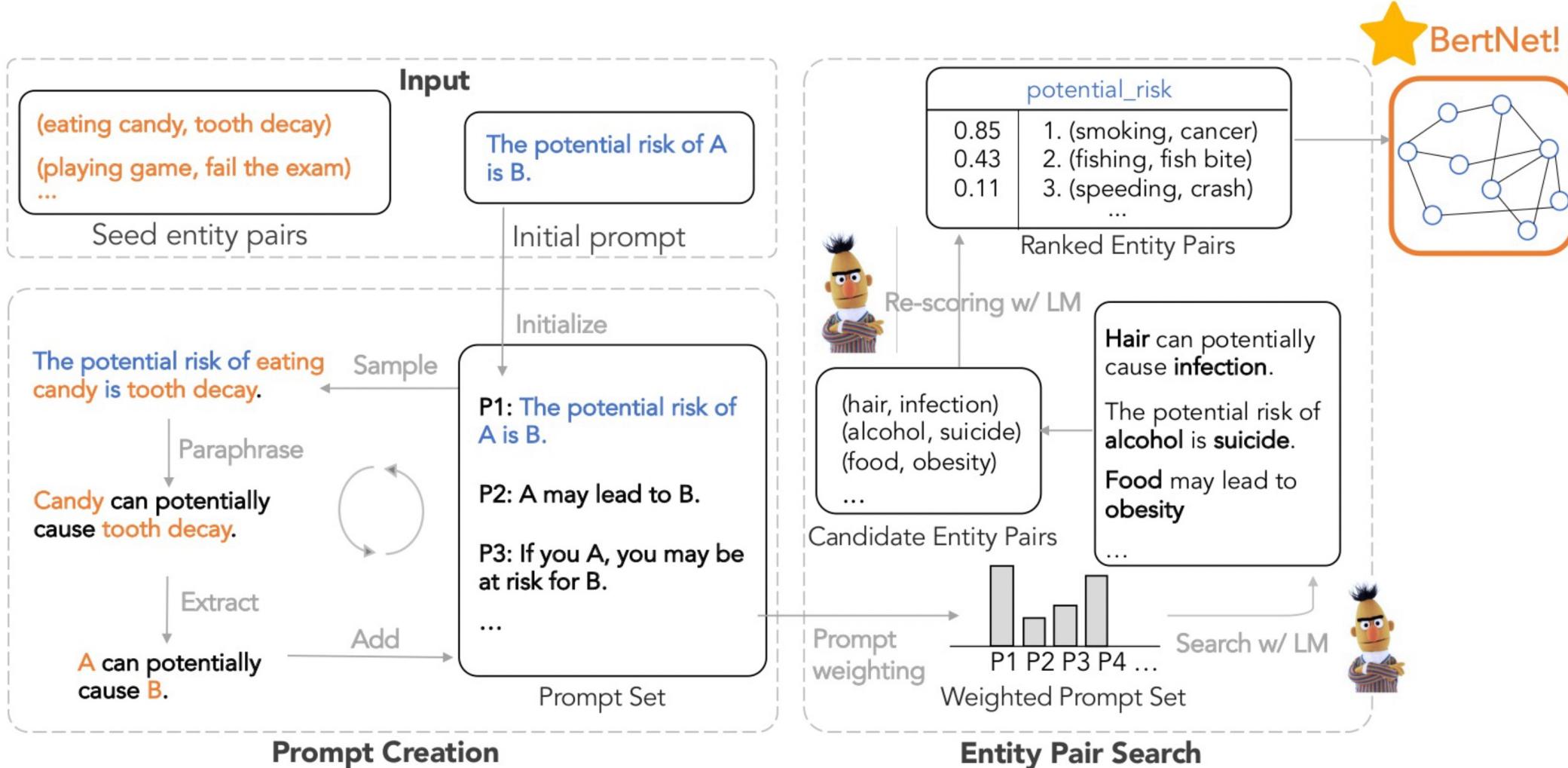
Fig 4. Example of generated depth-2 knowledge graph for entity Alan Turing

Limitations:

Error propagation:

Once wrong entity is predicted, the number of mistakes might increase

BERTNet



Research Questions

Q1. Can German clinical knowledge be extracted from pre-trained language models?

Q2. What is the level of accuracy achieved by the resulting Knowledge Graph, and which data and models exhibit the highest performance?

Q3. What is the best evaluation strategy?

- Resources:
 - Data, Models, Domain adaption

Resources - Models

- medBERT.de
- gBERT-large
- gBERT-base

Dataset	Size
OSCAR	145
OPUS	10
Wikipedia	6
OpenLegalData	2.4

Table 1a. Training data with size in GB.

	Params	GermEval18 (Coarse)	GermEval18 (Fine)	GermEval14	Averaged F1
DBMDZ BERT _{Base}	110m	75.23	47.39	87.90	70.17
deepset BERT _{Base}	110m	74.7	48.8	86.87	70.12
mBERT _{Base}	172m	70.00	45.20	87.44	67.55
XLM-Roberta _{Large}	550m	78.38	54.1	87.07	73.18
GBERT _{Data}	110m	74.51	48.01	87.41	69.97
GBERT _{WWM}	110m	76.48	49.99	87.80	71.42
GBERT _{Data + WWM}	110m	78.17	50.90	87.98	72.35
GBERT _{Large}	335m	80.08	52.48	88.16	73.57
GELECTRA	110m	76.02	42.22	86.02	68.09
GELECTRA _{Data}	110m	76.59	46.28	86.02	69.63
GELECTRA _{Large}	335m	80.70	55.16	88.95	74.94
Previous SoTA		76.77 (TU Wien)	52.71 (uhhLT)	84.65 (FLAIR)	

Table 1b. Results for gBERT-large und gBERT-base

medBERT.de:

A Comprehensive German BERT Model for the Medical Domain

- Trained on a diverse set of medical texts:
 - Scientific texts, medical books, and hospital data from various medical domains

Source	No. Documents	No. Sentences	No. Words	Size (MB)
DocCheck Flexikon	63,840	720,404	12,299,257	92
GGPONC 1.0 [2]	4,369	66,256	1,194,345	10
Webcrawl [24]	11,322	635,806	9,323,774	65
PubMed abstracts	12,139	108,936	1,983,752	16
Radiology reports	3,657,801	60,839,123	520,717,615	4,195
Spinger Nature	257,999	14,183,396	259,284,884	1,986
Electronic health records [22]	373,421	4,603,461	69,639,020	440
Doctoral theses	7,486	4,665,850	90,380,880	648
Thieme Publishing Group	330,994	10,445,580	186,200,935	2,898
Wikipedia	3,639	161,714	2,799,787	22
Summary	4,723,010	96,430,526	1,153,824,249	10,372

Table 2. Data Sources for medBERT.de^[1]

Model	AUROC	Macro F1	Micro F1	Precision	Recall
Chest CT					
GottBERT	92.48	69.06	83.98	76.55	65.92
BioGottBERT	92.71	69.42	83.41	80.67	65.52
Multilingual BERT	91.90	66.31	80.86	68.37	65.82
German-MedBERT	92.48	66.40	81.41	72.77	62.37
<i>medBERT.de</i>	96.69	81.46	89.39	87.88	78.77
<i>medBERT.de_{dedup}</i>	96.39	78.77	89.24	84.29	76.01
Chest X-Ray					
GottBERT	83.18	64.86	74.18	59.67	78.87
BioGottBERT	83.48	64.18	74.87	59.04	78.90
Multilingual BERT	82.43	63.23	73.92	56.67	75.33
German-MedBERT	83.22	63.13	75.39	55.66	78.03
<i>medBERT.de</i>	84.65	67.06	76.20	60.44	83.08
<i>medBERT.de_{dedup}</i>	84.42	66.92	76.26	60.31	82.99
ICD-10 code classification on discharge notes					
GottBERT	77.23	18.32	51.23	38.30	14.27
BioGottBERT	78.01	17.96	50.56	35.97	13.95
Multilingual BERT	76.64	19.48	51.19	38.39	15.60
German-MedBERT	75.44	23.41	53.63	41.39	18.94
<i>medBERT.de</i>	80.78	23.41	53.84	41.42	18.75
<i>medBERT.de_{dedup}</i>	80.84	21.44	52.46	40.45	17.04

Table 3. Results for medBERT.de^[1]

German clinical data

- Cardiology Doctoral Letters:
 - CARDIO:DE – 500 cardio discharge letters
 - MieDEEP – 500 cardio discharge letters
- Cardiology guidelines:
 - DGK: Leitlinien der Deutschen Gesellschaft für Kardiologie
 - 128 guidelines
- Oncology guidelines:
 - GGPONC 2.0 – 10,193 reports
- Oncology reports:
 - BRONCO150- 150 cancer patient reports

	Guideline	Year	Files
1	• Pancreatic cancer	2013	292
2	• Penis cancer	2020	167
3	• Psycho-oncology	2014	121
4	○ Oral cavity cancer	2021	132
5	• Malignant ovarian tumors	2020	195
6	• Anal cancer	2020	216
7	• Chronic lymphocytic leukemia	2018	285
8	• Laryngeal cancer	2019	189
9	• Follicular lymphoma	2020	296
10	• Oesophageal cancer	2018	172
11	○ Hodgkin lymphoma	2020	253
12	○ Hepatocellular and biliary cancer	2021	263
13	• Testicular tumors	2020	315
14	• Prevention of cervix cancer	2020	302
15	○ Renal cell carcinoma	2020	293
16	• Endometrial cancer	2018	317
17	• Stomach cancer	2019	246
18	• Adult soft tissue sarcomas	2021	407
19	• Actinic keratosis	2020	193
20	○ Malignant melanoma	2020	297
21	○ Cervical cancer	2021	415
22	• Colorectal cancer	2019	546
23	○ Prostate cancer	2021	351
24	• Supportive therapy	2020	819
25	• Lung cancer	2018	665
26	○ Breast cancer	2021	685
27	○ Bladder cancer	2020	364
28	○ Prevention of skin cancer	2021	370
29	○ Palliative medicine	2020	700
30	• Complementary medicine	2021	327
Total			10,193

	Iteration			
	1a	1b	2	3
Number of annotators	3	7	7	7
Number of documents	5	5	6	3
Number of sentences	149	149	158	67
Number of tokens	4206	4206	3725	1814
IAA (γ)	.75	.89	.93	.94
Specification	.71	.87	.91	.89
Finding	.82	.93	.95	.97
Diagnosis/Pathology	-	.91	.94	.96
Other Finding	-	.85	.87	.91
Substance	.92	.99	.98	.99
Clinical Drug	-	.97	.98	1.00
Nutrient/Body Subs.	-	.99	.99	.98
External Substance	-	.96	-	1.00
Procedure	.82	.93	.96	.96
Therapeutic	-	.95	.96	.96
Diagnostic	-	.89	.98	.93
IAA ($_{u\alpha}$)	.56	.71	.79	.85

Fig 5. Overview of guidelines in the current GGPONC release[1]

Frau
Dr. med. Paul Beispiel
Musterplatz 1
56789 Beispielstadt

Test-Klinik
Zentrum für Kardiologie

Klinik für Kardiologie
Station II

Dr. med. Muster
Ärztlicher Direktor
Station II

Station Sowieso
Beispielstr. 123
12345 Musterstadt
Tel +123 23 45 67
Fax +123 23 45 66

01.01.2010

CARDIO:DE annotations

Nachrichtlich:
Herrn Max Mustermann, Beispielplatz 1, 12345 Musterstadt

Sehr geehrter Herr Kollege Muster,

wir berichten über Ihre Patientin Frau Maxima Musterfrau geboren am 01.01.1970, wohnhaft in 12345 Musterstadt, Beispielstr. 1, die sich vom bis in unserer stationären Behandlung befand.

Diagnosen:

- Schwerer Infarkt der ... am 01.02
- Cvr: **Hyperlipidämie, Nikotinkonsum seit 01.01.1980, 30 py.**
- Allergien: Hausstaub

Anamnese:

Die stationäre Übernahme von Frau Musterfrau erfolgte über die Chirurgie. Die Patientin klagt über **Tachykardien**. Auf gezielte Nachfrage eingeschränkte Belastbarkeit, **belastungsabhängiges thorakales Druck- und Engegefühl** außerdem **progrediente Belastungsdyspnoe**. Es bestehen Ödeme bds., kein Schwindelgefühl, keine Synkopen.

Wir danken für die vertrauliche Zusammenarbeit und stehen bei Rückfragen selbstverständlich jederzeit gerne zur Verfügung.

Labor:

Bezeichnung	Wert	Datum
Abc	123	01.01.2010

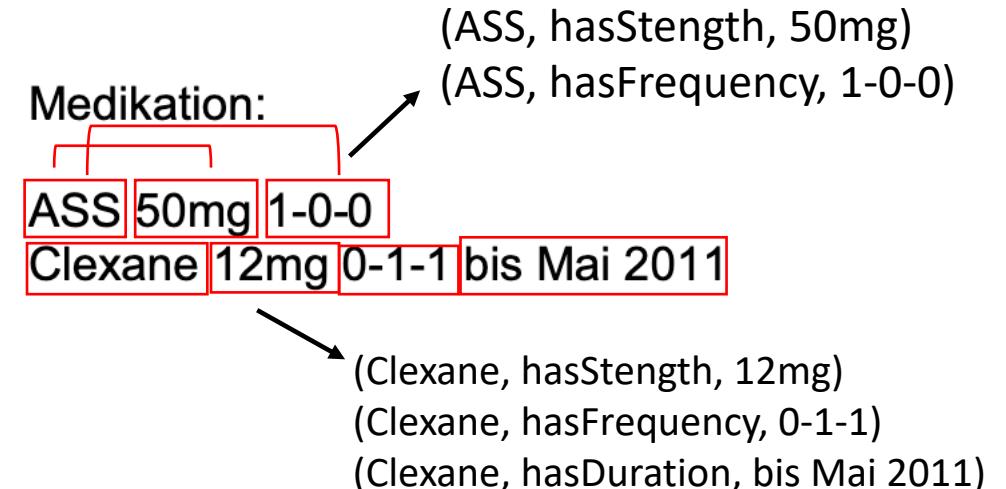
Medikation:

ASS 50mg 1-0-0
Clexane 12mg 0-1-1 bis Mai 2011

Mit freundlichen Grüßen

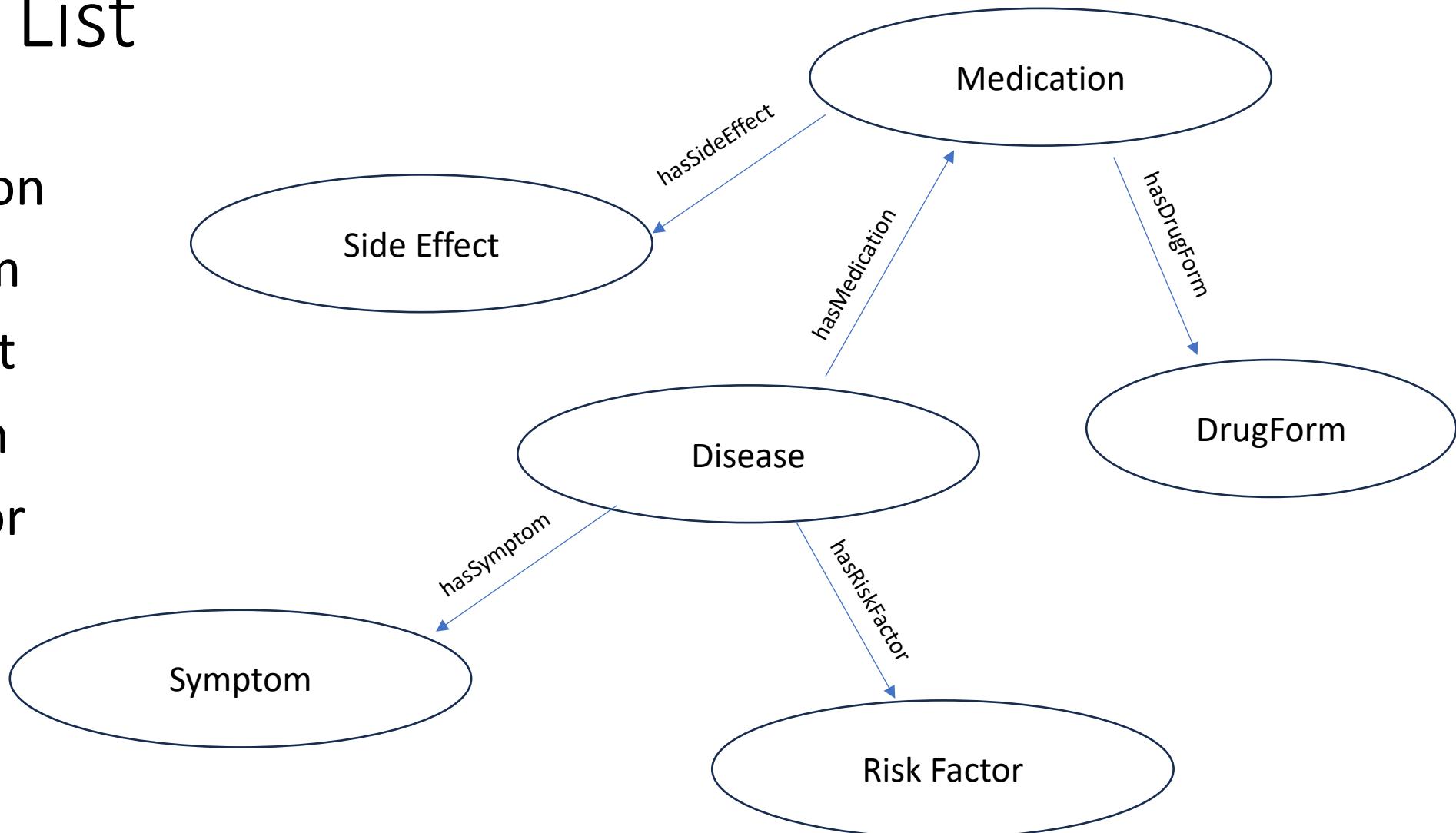
Medication information:

(Activeing, Dosage, Drug, Duration, Form, Frequency, Reason, Route, Strength)



Relations List

- hasMedication
- HasDrugForm
- HasSideEffect
- HasSymptom
- HasRiskFactor



Experiments

- Vanilla vs further pre-trained models
- Contextualized vs non-contextualized prompt
- Contextualized vs few-shot-prompt

Prompt Samples

1. Das <ENT1> wird bei der Krankheit <ENT0> angewendet.
2. Das in der Kardiologie verwendete Medikament <ENT1> dient der Behandlung von Krankheit <ENT0> .
3. Medikament: Warfarin, Krankheit: Vorhofflimmern. Das Medikament <ENT1> wird zur Behandlung der Krankheit <ENT0> in der Kardiologie eingesetzt.
4. In der Kardiologie wird das Medikament Warfarin bei Vorhofflimmern eingesetzt. Auch das Medikament Amiodaron wird zur Behandlung der Arrhythmie verwendet. Das Medikament Aspirin wird bei Arteriosklerose angewendet. Das Medikament <ENT1> wird zur Behandlung der Krankheit <ENT0> in der Kardiologie eingesetzt.

gBERT-base vs gBERT-large vs medBERT.de

Ours (Top 20)
['athleten', 'Doping']
['sauerstoff', 'Bakterien']
['Cannabis', 'Krankheiten']
['Cannabis', 'Krebs']
['Cannabis', 'Schmerzen']
['Cannabis', 'Patienten']
['Alkohol', 'Krebs']
['Alkohol', 'Krankheiten']
['Alkohol', 'Schmerzen']
['Fleisch', 'haut']
['Cannabis', 'Beschwerden']
['Blut', 'Menschen']
['Wasser', 'Erkrankungen']
['Alkohol', 'Erkrankungen']
['Gold', 'Menschen']
['Zucker', 'Diabetes']
['Sauerstoff', 'Bakterien']
['Fleisch', 'Erkrankungen']
['Energie', 'die umwelt']
['Sauerstoff', 'Beschwerden']

gBERT-base

Ours (Top 20)
['Magen Darm', 'Magen Pack']
['Krebs', 'Sar Ag']
['Krebs', 'Sar Cor']
['Grippe', 'Spr Comp']
['Grippe', 'Spr Cor']
['Schmerzen', 'Eisen']
['Schmerzen', 'Carmen']
['Nieren', 'Sar Vit']
['Schmerzen', 'Fieber']
['Schmerzen', 'Alkohol']
['Beschwerden', 'Nit Walk']
['Rauch Bedingt', 'Rauch Frei']
['Fieber', 'Eisen']
['Herz', 'Tore Ag']
['Schlaf', 'Nova Cor']
['Blut', 'Sot Cell']
['Erkrankungen', 'Eisen']
['Schlaf', 'Nova Plus']
['Krebs', 'Sar Plan']
['Leiden', 'Tore Ag']

gBERT-large

Ours (Top 20)
['Vorhofflimmern', 'Clopidogrel']
['Herzrhythmusstörungen', 'Magnesium']
['Vorhofflimmern', 'Magnesium']
['Herzrhythmusstörungen', 'Clopidogrel']
['Herzinsuffizienz', 'Clopidogrel']
['Herzinsuffizienz', 'Magnesium']
['Vorhofflimmern', 'Ibuprofen']
['Hypertonie', 'Clopidogrel']
['Vorhofflimmern', 'Heparin']
['Herzrhythmusstörungen', 'Ibuprofen']
['Arteriosklerose', 'Clopidogrel']
['Herzinsuffizienz', 'Ibuprofen']
['Herzrhythmusstörungen', 'Heparin']
['Herzinsuffizienz', 'Heparin']
['Hypertonie', 'Magnesium']
['Arteriosklerose', 'Ibuprofen']
['Hypertonie', 'Metformin']
['Arteriosklerose', 'Magnesium']
['Arteriosklerose', 'Heparin']
['Bluthochdruck', 'Vitamin Plus']

medBERT.de

Non- Contextualized vs Contextualized prompts

<ENT0> is recommended for a disease <ENT1>

Ours (Top 20)
['Magnesium', 'Angina Pectoris']
['Psychotherapie', 'Beschwerden']
['Metformin', 'kardiovaskulären risikofaktoren']
['Ibuprofen', 'Beschwerden']
['Sauerstoff', 'Atemnot']
['Psychotherapie', 'Schmerzen']
['aktuell methotrexat', 'Lungenfibrose']
['Cannabis', 'Beschwerden']
['Clopidogrel', 'Vorhofflimmern']
['Psychotherapie', 'Depression']
['Nikotin', 'kardiovaskulären risikofaktoren']
['Clopidogrel', 'Beschwerden']
['Prednisolon', 'akuten nierenversagen']
['Cannabis', 'Schmerzen']
['Sauerstoff', 'Dyspnoe']
['Physiotherapie', 'Rückenschmerzen']
['Paracetamol', 'Übelkeit']
['Ibuprofen', 'Atemnot']
['psychotherapeutische unterstützung', 'Depressionen']
['Psychotherapie', 'Schlafstörungen']

Non-contextualized prompt

In Cardiology, <ENT0> is recommended for a diease <ENT1>

Ours (Top 20)
['Magnesium', 'Herzrhythmusstörungen']
['Clopidogrel', 'Vorhofflimmern']
['Clopidogrel', 'Herzinsuffizienz']
['Kontrastmittel', 'Herzrhythmusstörungen']
['Metformin', 'Herzinsuffizienz']
['Kontrastmittel', 'Komplikationen']
['Magnesium', 'Herzinsuffizienz']
['Röntgen Thorax', 'Rundherde']
['Clopidogrel', 'Monate']
['Röntgen Thorax', 'Pneumothorax']
['Adenosin', 'weitere therapieansätze']
['Kontrastmittel', 'Beschwerden']
['Clopidogrel', 'Patienten']
['Metformin', 'Herzrhythmusstörungen']
['Cholesterin', 'kardiovaskuläre risikofaktoren']
['Adenosin', 'Beschwerden']
['Methotrexat', 'Lungenfibrose']
['Sauerstoff', 'Herzrhythmusstörungen']
['Kontrastmittel', 'Patienten']
['Methotrexat', 'Cyclophosphamid']

Contextualized prompt

Vanilla medBERT.de vs further pre-trained

Ours (Top 20)
['Arteriosklerose', 'Heparin']
['Arteriosklerose', 'Clopidogrel']
['Arteriosklerose', 'Statin']
['Arteriosklerose', 'Nikotin']
['Herzinsuffizienz', 'Heparin']
['Herzinfarkt', 'Heparin']
['Vorhofflimmern', 'Heparin']
['Bluthochdruck', 'Heparin']
['Herzinsuffizienz', 'Statin']
['Vorhofflimmern', 'Clopidogrel']
['Herzinfarkt', 'Clopidogrel']
['Herzinsuffizienz', 'Clopidogrel']
['Hypertonie', 'Mass Therapy']
['Herzinfarkt', 'Statin']
['Diabetes', 'Insulin']
['Vorhofflimmern', 'Statin']
['Bluthochdruck', 'Statin']
['Bluthochdruck', 'Clopidogrel']
['Arthrose', 'Ibuprofen']
['Hypertonie', 'Rot Gelb']

2 / top 5 correct
6 / top 10 correct
11 / top 20 correct

Ours (Top 20)
['Vorhofflimmern', 'Clopidogrel']
['Herzrhythmusstörungen', 'Magnesium']
['Vorhofflimmern', 'Magnesium']
['Herzrhythmusstörungen', 'Clopidogrel']
['Herzinsuffizienz', 'Clopidogrel']
['Herzinsuffizienz', 'Magnesium']
['Vorhofflimmern', 'Ibuprofen']
['Hypertonie', 'Clopidogrel']
['Vorhofflimmern', 'Heparin']
['Herzrhythmusstörungen', 'Ibuprofen']
['Arteriosklerose', 'Clopidogrel']
['Herzinsuffizienz', 'Ibuprofen']
['Herzrhythmusstörungen', 'Heparin']
['Herzinsuffizienz', 'Heparin']
['Hypertonie', 'Magnesium']
['Arteriosklerose', 'Ibuprofen']
['Hypertonie', 'Metformin']
['Arteriosklerose', 'Magnesium']
['Arteriosklerose', 'Heparin']
['Bluthochdruck', 'Vitamin Plus']

4 / top 5 correct
7 / top 10 correct
11 / top 20 correct

hasMedication – Fine-tuned medBERT.de

“Citalopram 20 mg”



Die Stärke von Citalopram beträgt 20 mg. Das Medikament Citalopram wird zur Behandlung von Depressionen angewendet. Müdigkeit ist eine Nebenwirkung des Medikaments Citalopram. Die Darreichungsform für Citalopram sind Tabletten.

3 / top 5

7 / top 10

11 / top 20

Ours (Top 20)

```
[ 'Bluthochdruck', 'Ibuprofen' ]
[ 'Diabetes Mellitus', 'Metformin' ]
[ 'Bluthochdruck', 'Clopidogrel' ]
[ 'Vorhofflimmern', 'Clopidogrel' ]
[ 'Herzinsuffizienz', 'Ibuprofen' ]
[ 'Vorhofflimmern', 'Ibuprofen' ]
[ 'Bluthochdruck', 'Methotrexat' ]
[ 'Hypertonie', 'Ibuprofen' ]
[ 'Bluthochdruck', 'Magnesium' ]
[ 'Herzrhythmusstörungen', 'Magnesium' ]
[ 'Vorhofflimmern', 'Magnesium' ]
[ 'Herzinsuffizienz', 'Clopidogrel' ]
[ 'Hypertonie', 'Clopidogrel' ]
[ 'Herzrhythmusstörungen', 'Ibuprofen' ]
[ 'Herzinsuffizienz', 'Magnesium' ]
[ 'Herzrhythmusstörungen', 'Clopidogrel' ]
[ 'Hypertonie', 'Magnesium' ]
[ 'Vorhofflimmern', 'Methotrexat' ]
[ 'Herzrhythmusstörungen', 'Vitamin Beispiele' ]
[ 'Herzinfarkt', 'Vitamin Beispiele' ]
```

HasRiskFactor

```
Ours (Top 20)

['Herzinsuffizienz', 'Diabetes']
['Herzinsuffizienz', 'Übergewicht']
    ['Herzinfarkt', 'Diabetes']
['Herzinsuffizienz', 'Adipositas']
['Herzinsuffizienz', 'Rauchen']
    ['Herzinfarkt', 'Übergewicht']
['Herzinsuffizienz', 'Bluthochdruck']
    ['Bluthochdruck', 'Diabetes']
    ['Herzinfarkt', 'Adipositas']
['Herzinfarkt', 'chronische niereninsuffizienz']
    ['Herzinfarkt', 'Rauchen']
    ['Bluthochdruck', 'Übergewicht']
    ['Diabetes', 'Übergewicht']
    ['Diabetes', 'Adipositas']
['Herzrhythmusstörungen', 'Übergewicht']
    ['Bluthochdruck', 'Adipositas']
    ['Bluthochdruck', 'Rauchen']
['Herzrhythmusstörungen', 'Adipositas']
    ['Herzrhythmusstörungen', 'Rauchen']
['Schlaganfall', 'chronische niereninsuffizienz']
```

20 / top 20 correct

HasSymptom

```
+-----+
|   Ours (Top 20)
+-----+
|   ['Atemnot', 'Herzinfarkt']
|   ['Atemnot', 'Herzinsuffizienz']
|       ['Atemnot', 'Angina']
|   ['Dyspnoe', 'Herzinsuffizienz']
|       ['Atemnot', 'Krankheit']
|       ['Dyspnoe', 'Herzinfarkt']
|       ['Schwäche', 'Herzinfarkt']
|           ['Dyspnoe', 'Angina']
|   ['Schwäche', 'Herzinsuffizienz']
|       ['Schwäche', 'Angina']
|       ['Müdigkeit', 'Herzinfarkt']
|           ['Husten', 'Asthma']
|       ['Müdigkeit', 'Herzinsuffizienz']
|           ['Müdigkeit', 'Angina']
|           ['Husten', 'Angina']
|           ['Husten', 'Herzinfarkt']
|           ['Husten', 'Herzinsuffizienz']
|               ['Müdigkeit', 'Krankheit']
|               ['Bluthochdruck', 'Hypertonie']
|               ['Erschöpfung', 'Krankheit']
+-----+
```

5 / top 5
10 / top 10
17 / top 20

HasDrugForm

```
+-----+  
| Ours (Top 20) |  
+-----+  
| ['Ibuprofen', 'Tabletten'] |  
| ['Ibuprofen', 'Tropfen'] |  
| ['Paracetamol', 'Tabletten'] |  
| ['Magnesium', 'Tabletten'] |  
| ['Paracetamol', 'Tropfen'] |  
| ['Ibuprofen', 'Rezept'] |  
| ['Ciprofloxacin', 'Tabletten'] |  
| ['Ciprofloxacin', 'Tropfen'] |  
| ['Clopidogrel', 'frei verfügbar'] |  
| ['Magnesium', 'Tropfen'] |  
| ['Ibuprofen', 'Privat'] |  
| ['Ibuprofen', 'Morphin'] |  
| ['Schmerz', 'ganz überwiegend'] |  
| ['Paracetamol', 'Rezept'] |  
| ['Propofol', 'ganz überwiegend'] |  
| ['Magnesium', 'Rezept'] |  
| ['Paracetamol', 'arzneimittel'] |  
| ['Ciprofloxacin', 'Rezept'] |  
| ['Magnesium', 'natürlich wirkstoff'] |  
| ['Clopidogrel Steht', 'ab sofort'] |  
+-----+
```

5 / top 5
8 / top 10

HasSideEffect

```
+-----+  
| Ours (Top 20) |  
+-----+  
| ['Kalium', 'Schwindel'] |  
| ['Kalium', 'Herzrhythmusstörungen'] |  
| ['Furosemid', 'Schwindel'] |  
| ['Furosemid', 'Herzrhythmusstörungen'] |  
| ['Magnesium', 'Schwindel'] |  
| ['Magnesium', 'Herzrhythmusstörungen'] |  
| ['Kalium', 'Müdigkeit'] |  
| ['Kalium', 'Kopfschmerzen'] |  
| ['Herzinsuffizienz', 'Herzrhythmusstörungen'] |  
| ['Herzinsuffizienz', 'Schwindel'] |  
| ['Magnesium', 'Kopfschmerzen'] |  
| ['Magnesium', 'Müdigkeit'] |  
| ['Furosemid', 'Müdigkeit'] |  
| ['Kalium', 'Schlafstörungen'] |  
| ['Furosemid', 'Übelkeit'] |  
| ['Furosemid', 'Kopfschmerzen'] |  
| ['Magnesium', 'Übelkeit'] |  
| ['Metformin', 'Schwindel'] |  
| ['Herzinsuffizienz', 'Müdigkeit'] |  
| ['Metformin', 'Übelkeit'] |  
+-----+
```

9/20 correct
2/5 correct
3/10 correct

Research Questions 1 + 2

Q1. Can German clinical knowledge be extracted from pre-trained language models?

- Yes, it can be. Moreover, by using a set threshold of the top 5 or top 10 pairs(in some cases), accurate matches can be extracted.

Q2. What is the level of accuracy achieved by the resulting Knowledge Graph, and which data and models exhibit the highest performance?

- medBERT.de further pre-trained on the cardio data

Knowledge Graph Construction

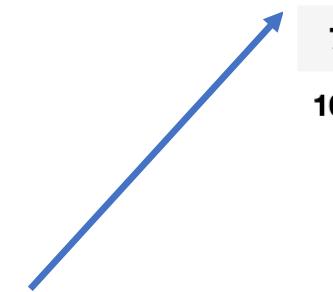
- Set a threshold to top 5 or even top 10 to ensure accurate relations.
- What if we had an extra filtering step?
 - Using GPT4 to remove wrong entities from the list
 - Keep the requested list short, to avoid misinformation from the gpt4

Filtered Results

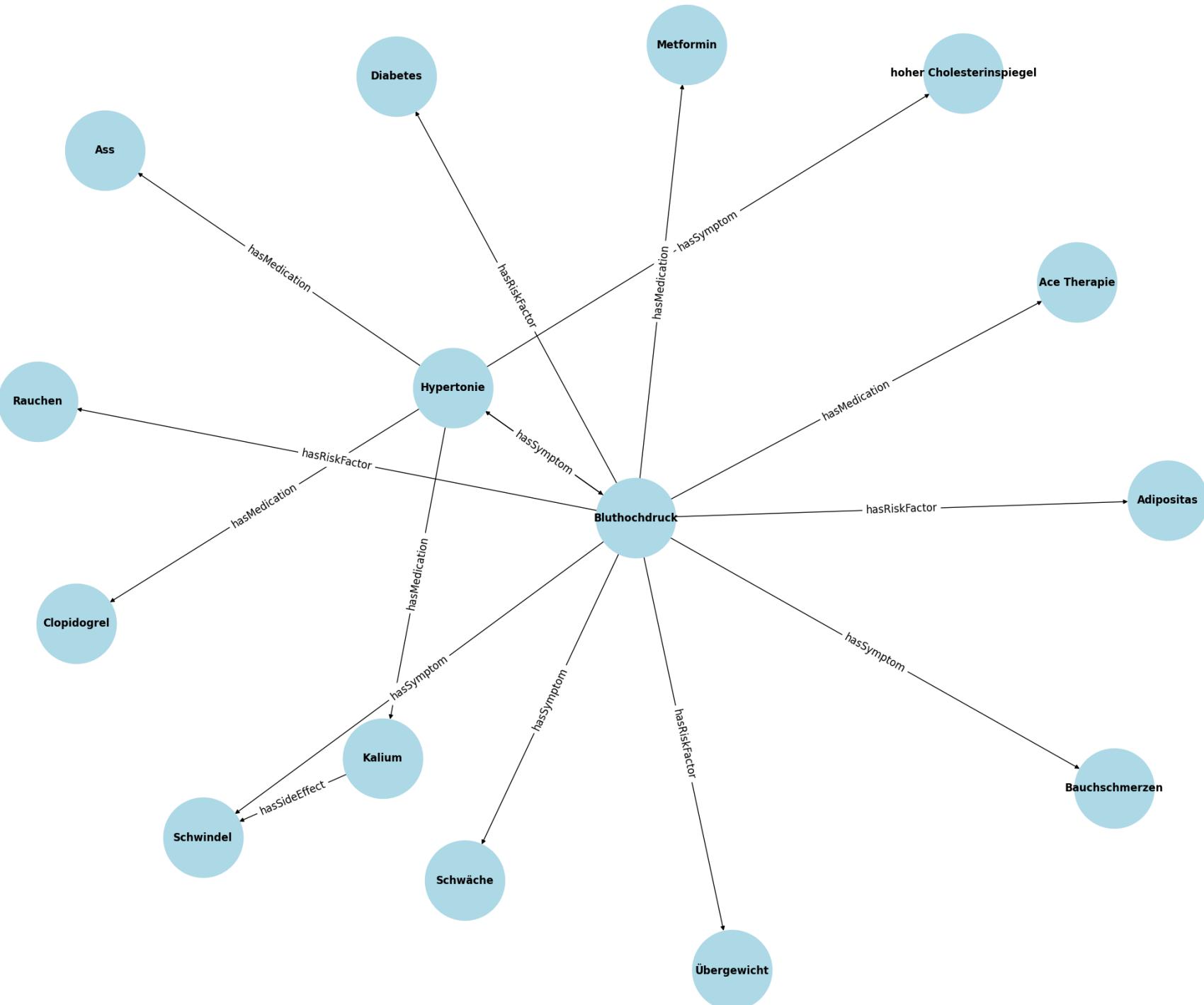
- HasMedication: 156 generated pairs --> 50 correct pairs
- HasDrugForm: 116 generated pairs -- > 15 correct pairs
- HasRiskFaktor: 68 generated pairs -- > 45 correct pairs
- HasSideEffect: 100 generated pairs -- > 18 correct pairs
- HasSymptom: 125 generated pairs -- > 68 correct pairs

The Resulting Database

	sourceid	relationship	destinationid
0	Vorhofflimmern	hasMedication	Clopidogrel
1	Herzrhythmusstörungen	hasMedication	Magnesium
2	Herzinsuffizienz	hasMedication	Clopidogrel
3	Herzinsuffizienz	hasMedication	Magnesium
4	Hypertonie	hasMedication	Clopidogrel
...
190	Azathioprin	hasSideEffect	Durchfall
191	Erythromycin	hasSideEffect	Blutbild Veränderungen
192	Cisplatin	hasSideEffect	Durchfall
193	Thiamazol	hasSideEffect	Magen Nebenwirkungen
194	Schildrüsen Insulin	hasSideEffect	Schildrüsen Gewichtszunahme



	sourceid	relationship	destinationid
4	Hypertonie	hasMedication	Clopidogrel
6	Hypertonie	hasMedication	Kalium
48	Hypertonie	hasMedication	Ass
76	Hypertonie	hasSymptom	Bluthochdruck
78	Bluthochdruck	hasSymptom	Hypertonie
103	Hypertonie	hasSymptom	hoher Cholesterinspiegel
	sourceid	relationship	destinationid
10	Bluthochdruck	hasMedication	Metformin
32	Bluthochdruck	hasMedication	Ace Therapie
78	Bluthochdruck	hasSymptom	Hypertonie
79	Bluthochdruck	hasSymptom	Schwäche
80	Bluthochdruck	hasSymptom	Schwindel
89	Bluthochdruck	hasSymptom	Bauchschmerzen
141	Bluthochdruck	hasRiskFactor	Diabetes
142	Bluthochdruck	hasRiskFactor	Übergewicht
143	Bluthochdruck	hasRiskFactor	Adipositas
144	Bluthochdruck	hasRiskFactor	Rauchen



Conclusion

- German clinical knowledge can be extracted from a small pre-trained medBERT.de model.
- Few shot prompts with 3 examples + contextualized prompt gives the best results for a further pre-trained medBERT.de on cardio data
- A threshold needs to be set to ensure accurate results
- GPT4 can be used as the last step in the pipeline to filter the results
- Github: <https://github.com/dieterich-lab/knowledge-graph-extraction-from-llms>