

# Medical Data Science - Übungsabgabe

Krankheitsbild: Schwere Pneumonie  
Schwerpunkt: Klinische Scores und ihre Validierung

Simon Bosse 613202  
Anna Göing 606171  
Jonas Jakob 640960  
Vyvy Tran Ngoc 623420

20. Mai 2025

## Inhaltsverzeichnis

<b>1. Einleitung</b>	<b>2</b>
<b>2. Setup</b>	<b>2</b>
2.1. Datenbankeinrichtung . . . . .	2
2.2. Erste Exploration des Datensatzes . . . . .	2
2.3. CITI-Kurs und MIMIC-IV Zugang . . . . .	3
<b>3. Datenexploration und Identifikation relevanter Parameter</b>	<b>4</b>
3.1. Recherche relevanter medizinischer Konzepte . . . . .	4
3.2. Analyse der Datenbankstruktur . . . . .	7
3.3. Mapping von Standardkonzepten zu MIMIC-IV-Bezeichnungen . . . . .	8
<b>A. Anhang</b>	<b>11</b>
A2.3 CITI-Zertifikate aller Gruppenmitglieder . . . . .	11
A3.3 Vollständige Mapping-Tabelle . . . . .	12

# 1. Einleitung

Die schwere Pneumonie ist eine Entzündung des Lungengewebes, die zu respiratorischer Insuffizienz führen kann und häufig eine mechanische Beatmung erfordert. Sie ist eine der häufigsten Ursachen für die Aufnahme auf die Intensivstation und war vor allem während der Corona-Pandemie ein oft diskutiertes Gesprächsthema in den Nachrichten, insbesondere in Bezug auf die Triage.

Im Rahmen der Übung könnte es interessant sein, im Hinblick auf die Priorisierung von Patienten und die grundsätzliche Entscheidung einer Aufnahme auf die Intensivstation, die verschiedenen Scoring-Systeme für Pneumonie zu betrachten und somit die Methoden für die Datenaufbereitung, sowie auch die darauf angewendete Berechnung der Scores zu betrachten.

- **Github Repository:** medicaldata (öffentlich zugänglich)

## 2. Setup

### 2.1. Datenbankeinrichtung

Die Software von Docker wurde für das entsprechende System heruntergeladen und installiert, sowie die Version mit dem `docker --version` Befehl getestet und ein PostgreSQL-Container konnte mithilfe der bereitgestellten `docker-compose.yml` gestartet werden. Das github repository wurde heruntergeladen und die benötigten Dateien `create.sql` und `load_gz.sql` zusammen mit dem Demo-MIMIC Datensatz in den Container geladen, so dass die Indizierung direkt vor Ort stattfinden konnte. Somit war der Zugriff auf den Datensatz eingerichtet.

Die Installation wurde mit dem Ausführen der folgenden Abfragen verifiziert:

- Anzahl der Patienten: 100 Patienten
- Anzahl der verschiedenen Diagnosen: 1472 Diagnosen

### 2.2. Erste Exploration des Datensatzes

Die Recherche nach ICD-9 und/oder ICD-10 Codes führte zu der folgenden Liste im Zusammenhang mit Pneumonie:

- **ICD-9:** Diseases of the respiratory system 460-519; Pneumonia und Influenza 480 - 488; Pneumonia 480 - 486
- **ICD-10:** Grippe und Pneumonie J09 - J18; Pneumonie J13 - J18

Die Exploration des Datensatzes ergab, dass exakt 8 verschiedene ICD-9 Codes für Pneumonie vorkommen, davon am häufigsten war 486 *Pneumonia, organism unspecified*, sowie zweimal 4846 *Pneumonia in aspergillosis*, während die anderen Codes nur jeweils einmal

vorkamen. Diese verschiedenen ICD-9 Codes beziehen sich auf exakt 11 Patienten. Die ICD-10 Codes gaben exakt zwei Ergebnisse für genau zwei Patienten. Während der Ausführung mussten die vorgegebenen SQL-Schemata jeweils um die dazugehörige Tabelle ergänzt werden, um die korrekte Relation aufzurufen. Beispielsweise *FROM mimic\_hosp.diagnoses\_icd* anstatt direkt *diagnoses\_icd* Liste der ICD-Codes nach der Datenexploration im SQL-Listenschema:

- **ICD-9:** '4809', '48020', '48241', '48242', '4249', '4829', '4846', '486'
- **ICD-10:** 'J15211', J159'

Diese Listen werden während der folgenden Übungsaufgaben benutzt werden, um auf die spezifischen Patientenfälle mit Pneumonie zuzugreifen.

### 2.3. CITI-Kurs und MIMIC-IV Zugang

Das erste Modul des Kurses dreht sich um den *Belmont Report*, eine US-amerikanische Richtlinie zu ethischen Prinzipien für die Forschung am Menschen. Die drei zentralen ethischen Prinzipien sind hierbei: (1) Achtung der Menschenrechte, (2) Benefizienz (Risikominimierung und Vorteilsmaximierung) und (3) Gerechtigkeit. Das zweite Modul erläutert die Entstehungsgeschichte dieser und anderer ethischer Richtlinien für die Forschung am Menschen.

Darauffolgend wird auf IRBs (institutional review boards) eingegangen, deren Funktionsweise und Aufbau erklärt, und in welchen Fällen ein IRB eingesetzt werden muss. Das vierte Modul behandelt *record-based research*, es wird unter anderem erläutert, dass auch hier abgewogen werden muss zwischen persönlichen Interessen (wie Datenschutz und Vertraulichkeit) und dem Lösen von Forschungsfragen.

Im Folgenden werden Arbeitsweisen und Vorteile der Gen(om)forschung vorgestellt, aber auch damit verbundene ethische Probleme angesprochen. In Modul 6 werden verschiedene vulnerable Gruppen vorgestellt und ethische Implikationen (auch mit Hinsicht auf die Prinzipien des *Belmont Report*) diskutiert. Darauffolgend werden dann Regularien wie der HIPAA-Standard vorgestellt, die sensible Gesundheitsdaten schützen sollen.

Im achten Modul werden verschiedene Arten von Interessenskonflikten vorgestellt (z.B. monetäre) und wie diese reguliert werden können (z.B. COI committees).

Im folgenden ist der aktuelle Status der Absolvierung des CITI-Kurses und des Zugangs zum Datensatz zum Zeitpunkt der Abgabe dokumentiert:

	CITI-Kurs	MIMIC-IV Zugang	Datum Zugang
Simon Bosse	Absolviert	Vorhanden	29.04.2025
Anna Göing	Nicht absolviert	Nein	-
Jonas Jakob	Nicht absolviert	Nein	-
Vyvy Tran Ngoc	Absolviert	Abgelehnt	-

Tabelle 1: Übersicht über Zugänge zum MIMIC-IV Datensatz

### 3. Datenexploration und Identifikation relevanter Parameter

#### 3.1. Recherche relevanter medizinischer Konzepte

Die Suche nach relevanten, medizinischen Konzepten und deren standardisierten Bezeichnungen für schwere Pneumonie begann mit einer groben Suche auf *Pubmed*, in der nur nach dem Krankheitsbild gesucht wurde, ohne komplexere Verknüpfungen. Der Suchbegriff: „*Severe Pneumonia*“. Diese Suche lieferte 5,521 Ergebnisse. Da das manuelle Durchsuchen dieser Menge an Artikeln nicht zielführend ist, wurde die Suche weiter eingeschränkt.

Als nächstes haben wir nach Artikeln gesucht die zusätzlich zu dem Namen des Krankheitsbilds, klinische Parameter oder Marker erwähnen. Dafür wurde die Suche mit folgendem Suchbegriff erneut durchgeführt: „*Severe Pneumonia AND (biomarkers OR clinical parameters)*“. Diese Suche lieferte 343 Ergebnisse. Diese Artikel wurden nun grob manuell gescannt und potentiell relevante Parameter, oder Scores wurden rausgeschrieben. Nach dieser Suche konnten wir einige wichtige Parameter (wie z.B. Age, Gender, Respiratory Rate, ...) identifizieren.

Zudem wurden der Score *CURB-65* und der *Pneumonia Severity Index* in mehreren Artikeln erwähnt. Daher wurde eine weitere Suche durchgeführt um nur Artikel zu finden die diese Scores erwähnen. Der Suchbegriff: „*Severe Pneumonia AND („CURB-65 OR „Pneumonia Severity Index“)*“. Diese Suche liefert auf *Pubmed* 115 Ergebnisse die wieder manuell durchsucht wurden. Bei der manuellen Suche wurde sich darauf fokussiert, welche Parameter für die Berechnung der Beiden Scores relevant sind.

Die folgende Ergebnistabelle 2 umfasst die gefundenen Überschneidungen zwischen den im zweiten Schritt identifizierten Werten und den für die Berechnung der Scores relevanten Parameter, die somit die wichtigsten medizinischen Konzepte für schwere Pneumonie darstellen sollten.

CONCEPT_ ID	CONCEPT_ NA- ME	DOMAIN_ ID	VOCABULARY_ ID	CONCEPT_ CODE	Bezeichnung der Literatur	in Medizinische Relevanz
4265453	Age	Meas Value	SNOMED	397669002	Age	CURB-65 und PSI/PORT Score
4135376	Gender	Observation	SNOMED	263495000	Sex	PSI/PORT Score
45883663	Nursing Home	Answer	LOINC	LA27-8	Nursing Home Resident	PSI/PORT Score
45881443	Confusion	Means Value	LOINC	LA7530-4	Confusion, mental instability	CURB-65 und PSI/PORT Score
4313591	Respiratory Rate	Measurement	SNOMED	86290005	Respiratory or breathing rate	CURB-65 und PSI/PORT Score
4152194	Systolic blood pressure	Measurement	SNOMED	271649006	Systolic Blood Pressure	CURB-65 und PSI/PORT Score
36716470	Temperature	Observation	SNOMED	722490005	Temperature	PSI/PORT Score
4224504	Pulse	Measurement	SNOMED	8499008	Pulse	PSI/PORT Score.

CONCEPT_ID	CONCEPT_NAME	DOMAIN_ID	VOCABULARY_ID	CONCEPT_CODE	Bezeichnung in der Literatur	Medizinische Relevanz
4017361	Blood urea nitrogen measurement	Measurement	SNOMED	105011006	Blood Urea Nitrogen, BUN	CURB-65 und PSI/PORT Score
4097822	pH measurement, arterial	Measurement	SNOMED	27051004	pH, arterial pH	PSI/PORT Score
4097430	Sodium measurement	Measurement	SNOMED	25197003	Sodium	PSI/PORT Score
4149519	Glucose measurement	Measurement	SNOMED	36048009	Glucose	PSI/PORT Score
4151358	Hematocrit determination	Measurement	SNOMED	28317006	Hematocrit	PSI/PORT Score
4103460	Oxygen measurement, partial pressure, arterial	Measurement	SNOMED	25579001	Partial pressure of arterial oxygen	PSI/PORT Score
254061	Pleural effusion	Condition	SNOMED	60046008	Pleural effusion on x-ray	PSI/PORT Score

Tabelle 2: Gefundene Medizinische Konzepte für schwere Pneumonie

### 3.2. Analyse der Datenbankstruktur

Der MIMIC-IV Datensatz beherbergt eine größere Anzahl an Tabellen, von denen nur ein Teil in wirklicher Relevanz zum ausgewählten Krankheitsbild und den im vorherigen Abschnitt identifizierten Parametern steht.

In diesem Sinne begann die Datenexploration des Datensatzes mithilfe von *DBeaver* mit einem ersten Fokus auf die Ortung der spezifischen Patientendaten und einzigartige Identifier, um die Zusammenhänge im weiteren Verlauf der Analyse verstehen zu können. Dies lies sich recht schnell innerhalb der Tabelle *patients* unter dem Primärschlüssel *subject\_id* finden, welcher jedem Patienten eine einzigartige Identifikation zuweist. Dort ebenso vorhanden sind die Daten zum Geschlecht (*gender*) und anonymisierte Informationen zum Alter (*anchor\_age*), die für das gewählte Krankheitsbild von besonderem Interesse sind.

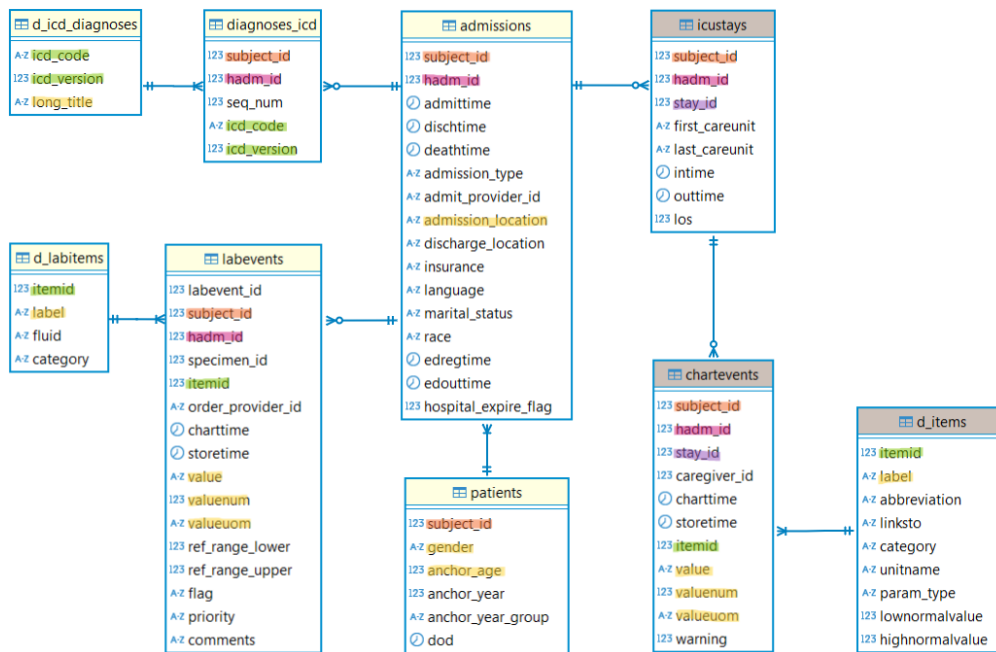


Abbildung 1: ER-Diagramm der identifizierten Tabellen des MIMIC-IV Datensatzes. Primär- und Fremdschlüssel wurden in entsprechenden Farben zwischen den Tabellen bunt markiert, wobei alle gelben Markierungen die Spaltennamen der relevanten Parameter für schwere Pneumonie sind.

Davon ausgehend geriet die Tabelle *admissions* in den Fokus, die für jede Ankunft im Krankenhaus der *subject\_id* eine Aufenthaltsid (*hadm\_id*) zuweist, die gemeinsam in Kombination alle weiteren Tabellen unter dem *mimiciv\_hosp* Schema als Schlüssel gelten. So lässt sich auch in der Tabelle *diagnoses\_icd* das zugehörige Krankheitsbild mithilfe der *d\_icd\_diagnoses* Tabelle identifizieren, wodurch der Datensatz auf die gewünschte Diagnose eingeschränkt werden kann. Ebenso lassen sich durch die Identifier aufgezeichnete Vital- und Laborwerte die während des Aufenthalts des Patienten im Krankenhaus ge-

messen wurden, innerhalb der Tabelle *labevents* finden, wo die *itemid* als Fremdschlüssel in der Tabelle *d\_labitems* auf das Label des Vorganges verweist. Dadurch sind die im Labor gemessenen Parameter in *labevents* unter *valuenum* und *valueuom* ablesbar und somit für die Analyse im Bezug auf das Krankheitsbild schwere Pneumonie aus diesem Datensatz nutzbar.

Ahnlich verhält es sich mit den Aufenthalten auf der Intensivstation unter dem Schema *mimici\_v\_icu*, die durch den Identifier des Patienten, der Krankenhaus-Aufenthalts ID und einer nun noch zusätzlich hinzukommenden *stay\_id* eindeutig zugewiesen werden können. Mithilfe dieser Dreierkombination enthält die Tabelle *chartevents* Einträge aller gemessenen Vital- und Laborwerten, die so wie in der Tabelle *labevents* abgelesen werden können und in der *d\_items* Tabelle durch den Primärschlüssel *itemid* das entsprechende Label zugewiesen bekommen. Dort lassen sich die relevanten Parameter für das ausgesuchte Krankheitsbild eingrenzen.

Die aus dieser Analyse der Datenbankstruktur identifizierten Tabellen und ihre Beziehung, sowie Zusammenhänge wurden im ER-Diagramm 1 dargestellt.

### 3.3. Mapping von Standardkonzepten zu MIMIC-IV-Bezeichnungen

Nachdem über die Analyse der Datenbankstruktur die wichtigen Tabellen identifiziert wurden, haben wir für alle Parameter SQL-Anfragen geschrieben, die dann für die relevanten Tabellen angepasst wurden. Die SQL-Anfragen folgen immer dem folgenden Schema:

```
SELECT <columns>
FROM <table_name>
WHERE <column_name> ILIKE '%-<keyword1>-%'
      OR <column_name> ILIKE '<keyword1>-%'
      OR <column_name> ILIKE '%-<keyword1>'
      OR <column_name> ILIKE '<keyword1>'
```

Für jede Anfrage wurden pro Keyword verschiedene mögliche Kombinationen mit Leerzeichen abgefragt. Pro identifiziertem Parameter wurde entschieden welche Synonyme in der ersten Version der Anfrage inkludiert werden sollten - diese orientieren sich nah an den Namen der Parameter in der SNOMED Datenbank.

Im ersten Schritt wurden alle Anfragen für unsere fünfzehn Parameter auf der Tabelle *d\_items* des ICU Schemas ausgeführt. In dem Fall, dass eine Anfrage keine Ergebnisse lieferte, es aber durchaus denkbar ist dass für diesen Parameter ein Ergebnis in der Tabelle vorliegen kann, wurde ein weiterer Rechenschritt durchgeführt, in welchem wir aus der Literatur relevante Synonyme identifiziert haben und die jeweilige SQL-Anfrage dann mit den weiteren Keywords angereichert haben. Dadurch sind teilweise sehr komplexe Anfragen entstanden, die dann aber auch für alle folgenden Tabellen weiter verwendet wurden.

Nachdem dann alle fünfzehn Anfragen für die *d\_items* Tabelle des ICU Schemas ausgeführt wurden, haben wir die Anfragen an die dictionary Tabellen des Hospital Schemas angepasst und dort ausgeführt. Vor allem war die Tabelle *labitems* relevant. Für



manche Parameter war es zudem sinnvoll in der Tabelle d-diagnoses eine Anfrage auszuführen. Dadurch dass die Anpassung der SQL Anfragen sehr aufwändig ist, wurden für die jeweiligen Parameter vor der Ausführung die Tabellen identifiziert, in denen es wahrscheinlich ist, auf Ergebnisse zu kommen.

Falls trotzdem kein Ergebnis in den dictionary Tabellen gefunden werden konnte, wurde ein weiterer Blick auf das Datenbankschema geworfen, gerade bei Informationen über den Patienten, zum Beispiel bezüglich eines Aufenthalts im Altenheim, lassen sich darüber die nötigen Informationen beschaffen, jedoch kann für diese Parameter dann leider kein Code identifiziert werden.

Des weiteren haben wir, falls diverse Varianten eines Ergebnisses gefunden wurden, alle Varianten in der Tabelle aufgeführt die wir als potentiell wichtig erachtet haben. Falls sich im späteren Verlauf des Projekts herausstellen sollte, dass einige Codes für Parameter komplett unnötig sind, können wir diese später immer noch streichen.

Das Resultat des Mappings der MIMIC Codes auf die SNOMED Codes, welche identifiziert wurden, findet sich im Appendix in Tabelle 3.

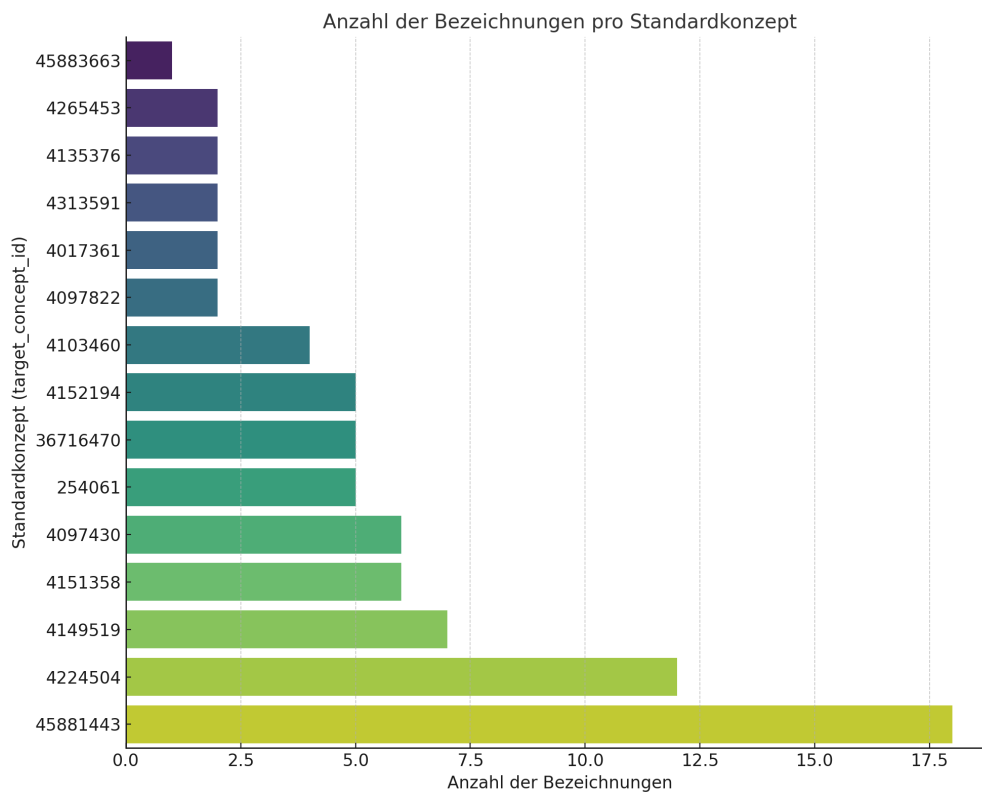


Abbildung 2: Balkendiagramm was das Mapping der Codes der Standardkonzepte auf die Zahl der jeweils identifizierten MIMIC Codes zeigt

Für die Visualisierung des erstellten Mappings wurde ein Balkendiagramm verwendet, da es 15 eins zu viele Beziehungen gibt, die visualisiert werden sollen. Da war das Balkendiagramm, abgebildet in Figur 2 die übersichtlichste Variante. Es ist anzumerken, dass bei den Konzepten, für die nur ein oder zwei Mappings gefunden wurden, die identifizierten Codes bei der weiteren Anwendung überprüft werden sollten, bzw. die benötigte Information ggf. anderswo in der Datenbank liegt.

## A. Anhang

### A2.3 CITI-Zertifikate aller Gruppenmitglieder



### A3.3 Vollständige Mapping-Tabelle

source_ code	source_code _description	source_ vocabula- ry_id	target_ con- cept_id	target_ vocabula- ry_id
none	Age	MIMIC	4265453	SNOMED
226984	Apache IV Age	MIMIC	4265453	SNOMED
none	Gender	MIMIC	4135376	SNOMED
226228	Gender	MIMIC	4135376	SNOMED
none	Nursing Home	MIMIC	45883663	LOINC
none	Confusion / Mental In- stability	MIMIC	45881443	LOINC
228395	Orientation to Place	MIMIC	45881443	LOINC
228394	Orientation to Person	MIMIC	45881443	LOINC
229381	Orientation	MIMIC	45881443	LOINC
223898	Orientation	MIMIC	45881443	LOINC
228396	Orientation to Time	MIMIC	45881443	LOINC
226104	Level of Consciousness	MIMIC	45881443	LOINC
229382	Orientation Score	MIMIC	45881443	LOINC
228688	Delirium	MIMIC	45881443	LOINC
2930	Delirium due to conditi- ons classified elsewhere	MIMIC	45881443	LOINC
2903	Senile dementia with delirium	MIMIC	45881443	LOINC
2931	Subacute delirium	MIMIC	45881443	LOINC
2982	Reactive confusion	MIMIC	45881443	LOINC
29281	Drug-induced delirium	MIMIC	45881443	LOINC
29011	Presenile dementia with delirium	MIMIC	45881443	LOINC
78097	Altered mental status	MIMIC	45881443	LOINC
29041	Vascular dementia, with delirium	MIMIC	45881443	LOINC
F05	Delirium due to known physiological condition	MIMIC	45881443	LOINC
none	Respiratory Rate	MIMIC	4313591	SNOMED
230040	Paradoxical breathing	MIMIC	4313591	SNOMED
225309	ART BP Systolic	MIMIC	4152194	SNOMED
227243	Manual Blood Pressure Systolic Right	MIMIC	4152194	SNOMED
220179	Non Invasive Blood Pressure Systolic	MIMIC	4152194	SNOMED
220050	Arterial Blood Pressure Systolic	MIMIC	4152194	SNOMED

source_ code	source_code _description	source_ vocabula- ry_id	target_ con- cept_id	target_ vocabula- ry_id
224167	Manual Blood Pressure Systolic Left	MIMIC	4152194	SNOMED
224027	Skin Temperature	MIMIC	36716470	SNOMED
223761	Temperature Fahren- heit	MIMIC	36716470	SNOMED
223762	Temperature Celsius	MIMIC	36716470	SNOMED
226329	Blood Temperature CCO (C)	MIMIC	36716470	SNOMED
50825	Temperature Blood Blood Gas	MIMIC	36716470	SNOMED
229770	Resting Pulse Rate (COWS)	MIMIC	4224504	SNOMED
223942	Graft/Flap Pulse	MIMIC	4224504	SNOMED
223936	Radial Pulse R	MIMIC	4224504	SNOMED
223948	Radial Pulse L	MIMIC	4224504	SNOMED
223941	Popliteal Pulse R	MIMIC	4224504	SNOMED
223946	Popliteal Pulse L	MIMIC	4224504	SNOMED
223949	Ulnar Pulse L	MIMIC	4224504	SNOMED
223945	Femoral Pulse L	MIMIC	4224504	SNOMED
223939	Brachial Pulse R	MIMIC	4224504	SNOMED
223944	Brachial Pulse L	MIMIC	4224504	SNOMED
223940	Femoral Pulse R	MIMIC	4224504	SNOMED
223938	Ulnar Pulse R	MIMIC	4224504	SNOMED
225624	BUN	MIMIC	4017361	SNOMED
51842	Bun	MIMIC	4017361	SNOMED
50820	pH	MIMIC	4097822	SNOMED
223830	PH (Arterial)	MIMIC	4097822	SNOMED
220645	Sodium (serum)	MIMIC	4097430	SNOMED
226534	Sodium (whole blood)	MIMIC	4097430	SNOMED
228389	Sodium (serum) (soft)	MIMIC	4097430	SNOMED
228390	Sodium (whole blood) (soft)	MIMIC	4097430	SNOMED
50983	Sodium	MIMIC	4097430	SNOMED
52623	Sodium	MIMIC	4097430	SNOMED
50809	Glucose	MIMIC	4149519	SNOMED
50931	Glucose	MIMIC	4149519	SNOMED
52569	Glucose	MIMIC	4149519	SNOMED
226537	Glucose (whole blood)	MIMIC	4149519	SNOMED
225664	Glucose finger stick (range 70-100)	MIMIC	4149519	SNOMED

source_ code	source_code _description	source_ vocabula- ry_id	target_ con- cept_id	target_ vocabula- ry_id
220621	Glucose (serum)	MIMIC	4149519	SNOMED
228388	Glucose (whole blood) (soft)	MIMIC	4149519	SNOMED
52028	Hematocrit Blood	MIMIC	4151358	SNOMED
51638	Hematocrit Blood	MIMIC	4151358	SNOMED
51639	Hematocrit Blood	MIMIC	4151358	SNOMED
51221	Hematocrit Blood	MIMIC	4151358	SNOMED
226540	Hematocrit (whole blood - calc)	MIMIC	4151358	SNOMED
220545	Hematocrit (serum)	MIMIC	4151358	SNOMED
220227	Arterial O2 Saturation	MIMIC	4103460	SNOMED
220277	O2 saturation pulseoxy- metry	MIMIC	4103460	SNOMED
223835	Inspired O2 Fraction	MIMIC	4103460	SNOMED
50817	Oxygen Saturation	MIMIC	4103460	SNOMED
51181	Malignant pleural effu- sion	MIMIC	254061	SNOMED
5119	Unspecified pleural effu- sion	MIMIC	254061	SNOMED
J910	Malignant pleural effu- sion	MIMIC	254061	SNOMED
J918	Pleural effusion in other conditions classified el- sewhere	MIMIC	254061	SNOMED
J91	Pleural effusion in con- ditions classified else- where	MIMIC	254061	SNOMED

Tabelle 3: Mapping der gefundenen Standardkonzepte auf die des Mimic Datensatzes