

Towards the integration of WordNet into ClinIDMap

Elena Zotova, Montse Cuadros and German Rigau
Vicomtech & University of Basque Country

Nombre

January 23-27, 2023
Donostia, Spain



Elena ZOTOVA

NLP Researcher, PhD Student
Tecnological Centre Vicomtech,
University of the Basque Country
ezotova@vicomtech.org



Dr Montse CUADROS

Senior NLP Researcher
Tecnological Centre Vicomtech

mcuadros@vicomtech.org



Dr German RIGAU

Deputy Director, Associate
Professor
HiTZ Basque Center for
Language Technologies,
University of the Basque Country
german.rigau@ehu.eus

Index

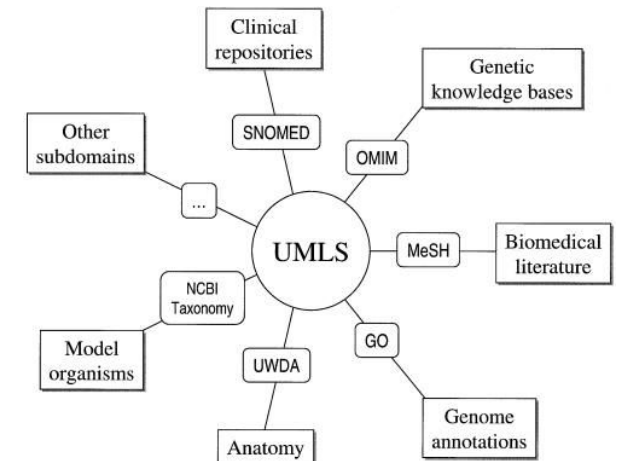
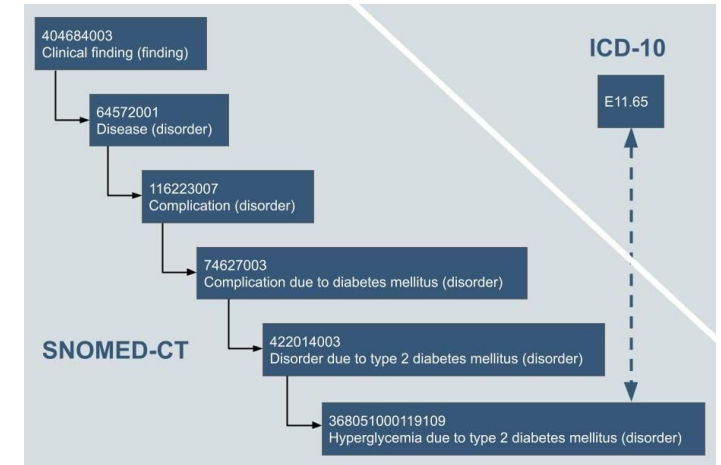
1. Introduction
2. Motivation
3. Background
4. Mapping method
5. Corpora Annotation
6. Conclusion

Index

1. Introduction
2. Motivation
3. Background
4. Mapping method
5. Corpora Annotation
6. Conclusion

Introduction

- Terminology processing is an important part of **clinical NLP**
- There are dictionaries, taxonomies and ontologies of different structure
- Similar concepts, synonyms, hierarchy
- Unique identifier - short description of a term / concept



- **ClinIDMap** - a tool for mapping identifiers between clinical ontologies and lexical resources (Zotova E et al, 2022)
- Interlinks identifiers from UMLS (CUI), SMOMED-CT, ICD-10 and Wikipedia/Wikidata items and WordNet
- The goal is to provide semantic interoperability across the clinical concepts from various knowledge bases
- Focus on Spanish language but the method is valid for any language
- Explore WordNet for clinical domain

Example from E3C Corpus

Durante los 5 años que permaneció en DP sufrió 10 **peritonitis** [C0031154], 8 por Staphilococcus aureus.

Translation:

During the 5 years on PD he suffered 10 **peritonitis** [C0031154], 8 of which were because of Staphylococcus aureus.

Magnini, B., Altuna, B., Lavelli, A., Speranza, M., & Zanolli, R. (2020). The E3C Project: Collection and Annotation of a Multilingual Corpus of Clinical Cases. Italian Conference on Computational Linguistics.

Example from E3C Corpus

Durante los 5 años que permaneció en DP sufrió 10 **peritonitis** [C0031154], 8 por Staphilococcus aureus.

Translation:

During the 5 years on PD he suffered 10 **peritonitis** [C0031154], 8 of which were because of Staphylococcus aureus.

The code **C0031154 (UMLS CUI)** corresponding to the Spanish term *peritonitis* can be mapped to:

- SNOMED CT code **235983003**;
- ICD-10-CM code **K65**;
- Wikipedia articles in **48 languages**;
- **synset 14376092-n** (WordNet 3.1);
- **synset 14352687-n** (WordNet 3.0).

Index

1. Introduction
2. Motivation
3. Background
4. Mapping method
5. Corpora Annotation
6. Conclusion

The task of clinical coding:

- Detect a span with clinical terminology (Sequence labelling model)
- Assign a high-level class: Diagnosis, Procedure, Anatomy etc.
- Link classified span with clinical taxonomy and assign an ID

Data for machine learning models in the clinical domain is especially difficult:

- Clinical information is private
- Manual annotation requires high-level expertise in medicine
- Few data is available for languages other than English
- If we have a corpus annotated with UMLS codes we transfer it to the corpus with ICD-10/SNOMED-CT/other codes

Index

1. Introduction
2. Motivation
3. Background
4. Mapping method
5. Corpora Annotation
6. Conclusion

Background

Concept alignment, or ontology alignment (also known as ontology matching): lexical matching, structural matching and logical reasoning.

Applications which use the resulting concept mapping to process biomedical text: MetaMap, I-MAGIC

I-MAGIC
Using 202109 release of the SNOMED CT to ICD-10-CM map

[About](#) | [Instructions](#) | [Demo](#)

The I-MAGIC (Interactive Map-Assisted Generation of ICD Codes) Algorithm utilizes the [SNOMED CT to ICD-10-CM Map](#) in a real-time, interactive manner to generate ICD-10-CM codes. This demo simulates a problem list interface in which the user enters problems with SNOMED CT terms, which are then used to derive ICD-10-CM codes using the Map.

Name: Gender: Date of Birth:

Mapping Problems to ICD-10-CM

SNOMED-CT	ICD-10-CM Code	ICD-10-CM Name	Optional refinement
Type 2 diabetes mellitus (44054006)	E11.9	Type 2 diabetes mellitus without complications	<input type="button" value="ICD notes"/>

Kin Wah Fung
[Lister Hill National Center for Biomedical Communications](#), [National Library of Medicine](#)
[NLM Web Policies](#), [NLM Support Center](#), [Accessibility](#)
[FOIA](#), [HHS Vulnerability Disclosure](#)
[NIH](#), [HHS](#), [USA.gov](#)

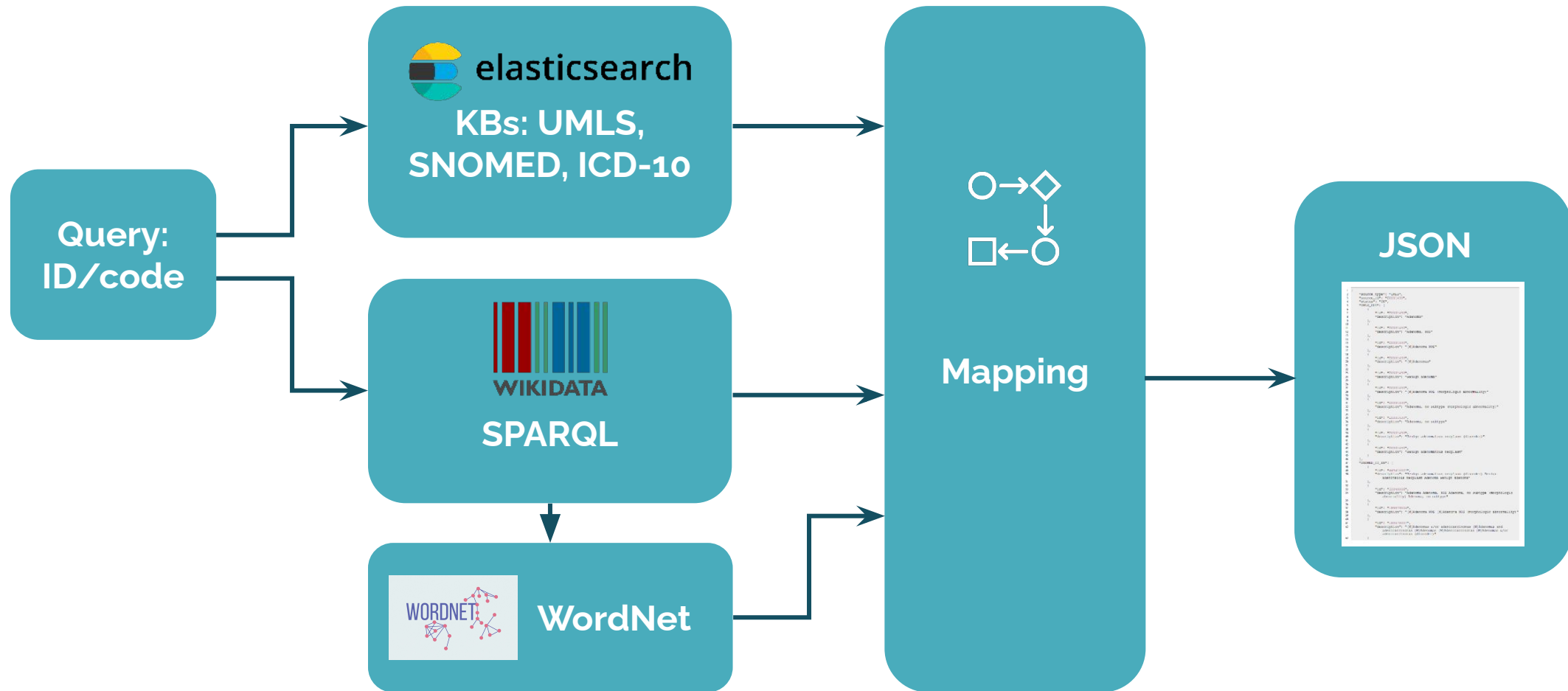
Attempts to use WordNets for clinical domain

- Medical WordNet had the goal of linking different terms, both professional terminology and general language. (Smith and Fellbaum, 2004)
- Explaining foreign clinical terms in Norwegian (Ingvaldsen and Veres, 2004)
- Disambiguation of UMLS mappings (Mougin et al., 2006).
- Aligning WordNet domains and Wikipedia categories to obtain domain corpora (Gella et al., 2014)

Index

1. Introduction
2. Motivation
3. Background
4. Mapping method
5. Corpora Annotation
6. Conclusion

High-level scheme



Mapping: Knowledge Bases

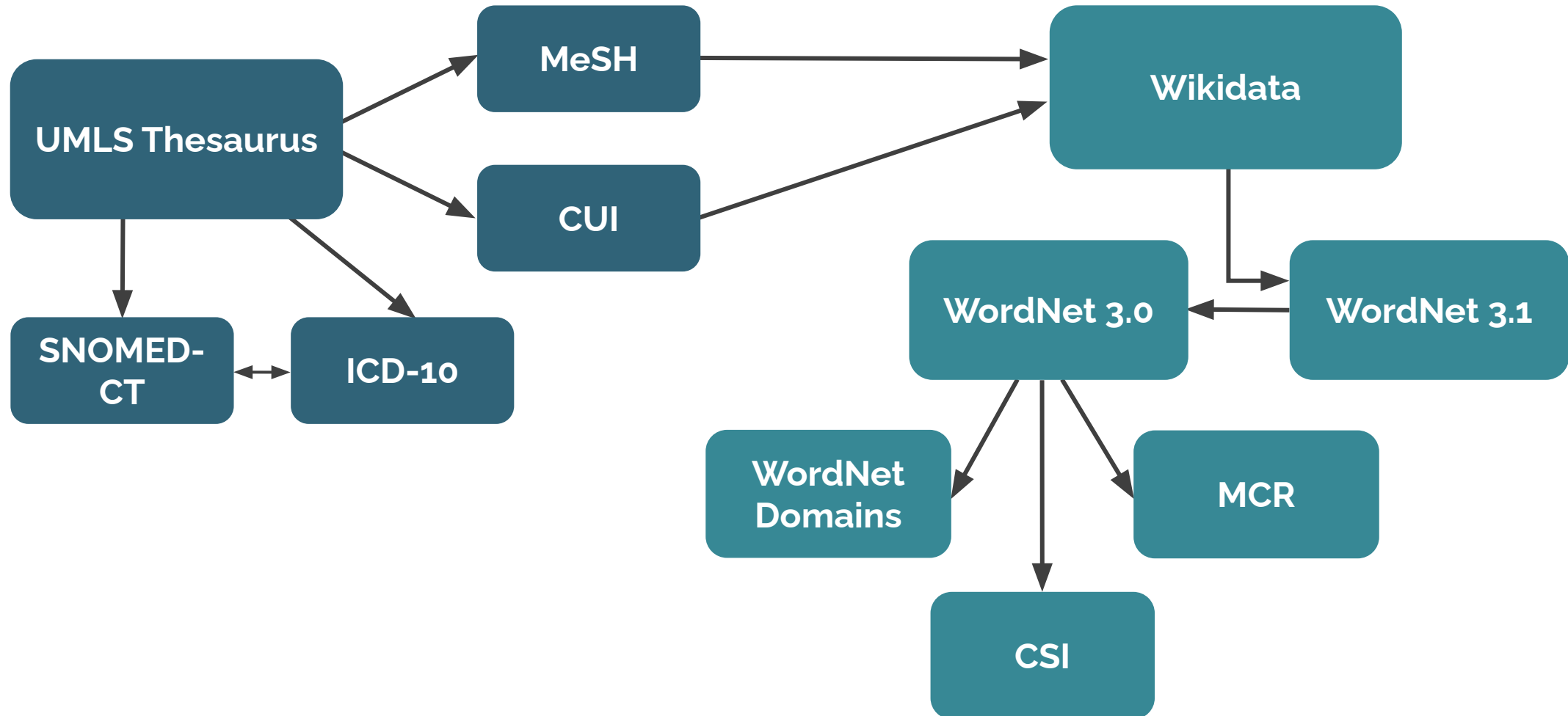


- UMLS Metathesaurus: mappings with different clinical ontologies
- SNOMED-CT (es, en)
- ICD-10-CM (es, en)
- ICD-10-PCS (es, en)
- SNOMED-CT to ICD-10 mapping (en)
- MeSH
- Wikidata
- Wikipedia
- WordNet 3.1 and 3.0

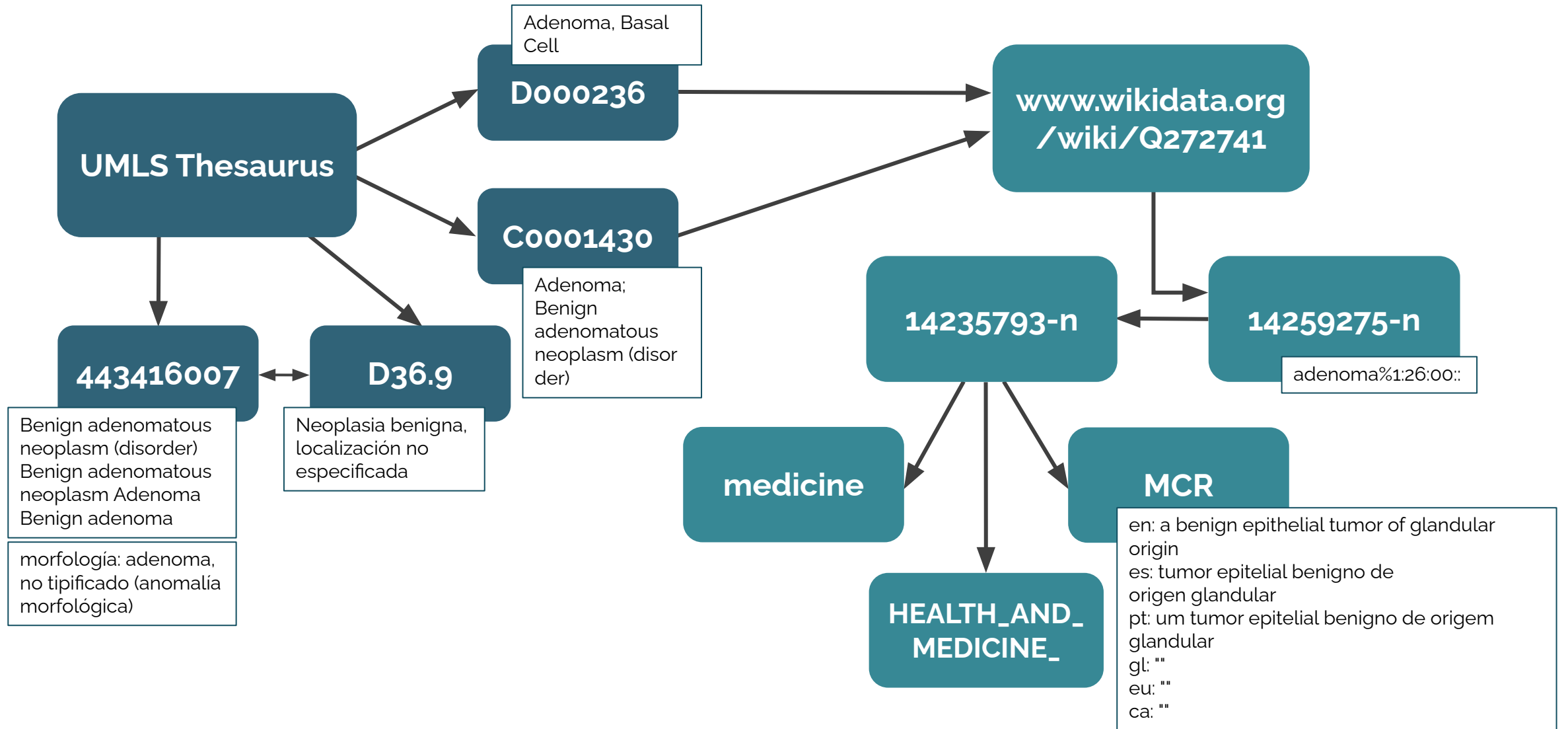
Mapping: Knowledge Bases

Knowledge base	Unique IDs
SNOMED-CT	489,141
ICD-10-CM	95,671
ICD-10-PCS	90,673
MeSH	347,565
UMLS CUI	1,106,486
SNOMED-CT 2 ICD-10	125,823
SNOMED-CT-ES 2 ICD-10	57,393
Wikidata	28,113
Wikidata+WordNet 3.1	26,953

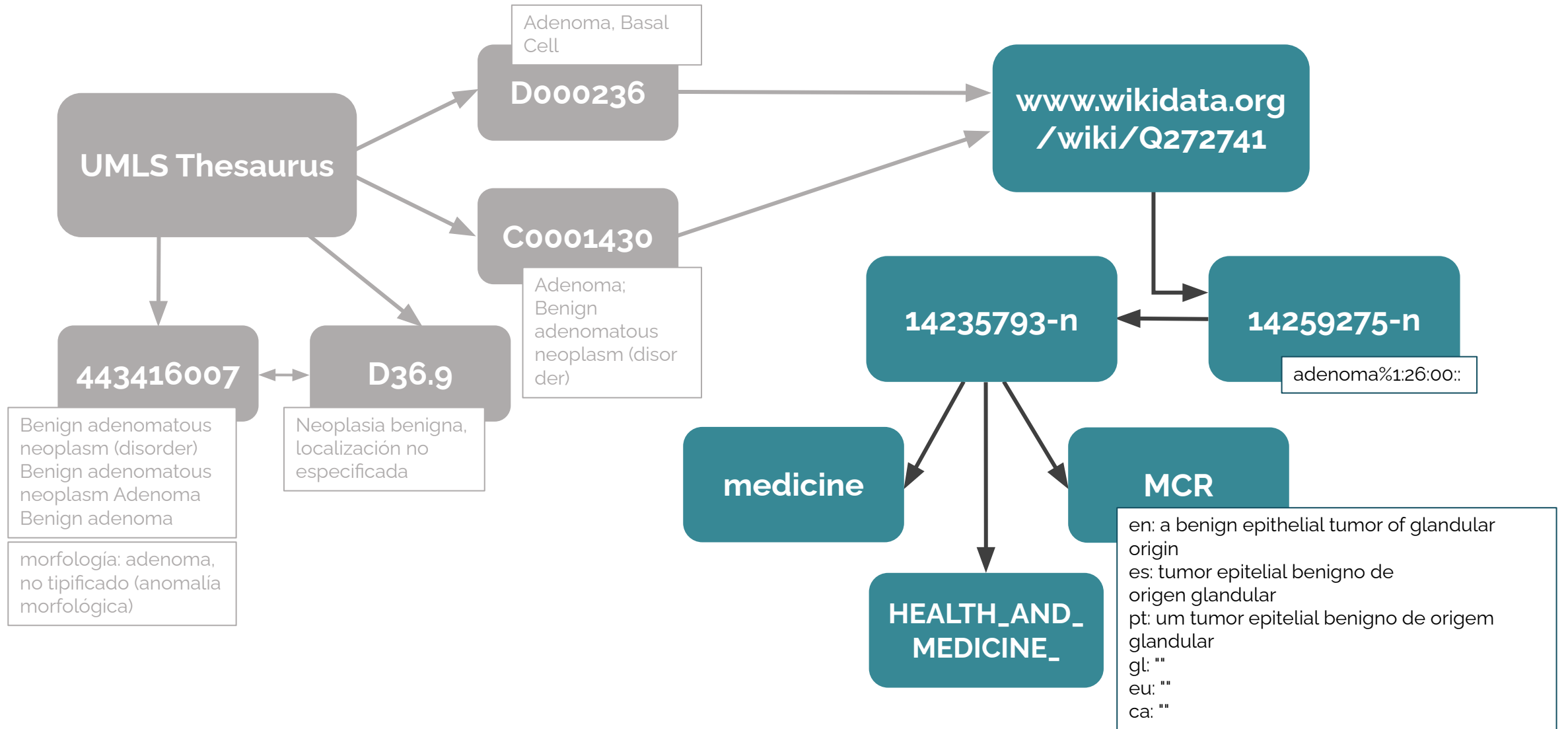
Mapping



Mapping



Mapping



Steps to integrate WordNets



1. Collect Wikidata items with SPARQL
2. Map WordNet 3.1. to WordNet 3.0 with index.sense
3. Map WordNet 3.0 to WN Domains and
4. Map to CSI (Coarse Sense Inventory)
5. Map to Multilingual WordNet (MCR)

Steps to integrate WordNets

Step 1. Collect Wikidata items with SPARQL

```
SELECT ?item ?itemLabel ?identificador_de_synset_de_WordNet_3_1 ?UMLS_CUI
?identificador_MeSH ?ICD_10 ?Snomed_CT WHERE {
  SERVICE wikibase:label { bd:serviceParam wikibase:language
"[AUTO_LANGUAGE],en". }
  OPTIONAL { ?item wdt:P8814 ?identificador_de_synset_de_WordNet_3_1. }
  OPTIONAL { ?item wdt:P2892 ?UMLS_CUI. }
  OPTIONAL { ?item wdt:P486 ?identificador_MeSH. }
  OPTIONAL { ?item wdt:P494 ?ICD_10. }
  OPTIONAL { ?item wdt:P4229 ?ICD_10_CM. }
  OPTIONAL { ?item wdt:P1690 ?ICD_10_PCS. }
  OPTIONAL { ?item wdt:P5806 ?Snomed_CT. }
}
```

Steps to integrate WordNets

Step 2. Map WordNet 3.1. to WordNet 3.0 with index.sense

WN 3.0 adenoma%1:26:00:: 14235793

WN 3.1. adenoma%1:26:00:: 14259275

Steps to integrate WordNet

Step 3. Domain Labels.

170 domain labels, clinical domain:

medicine, anatomy, pharmacy, health, biochemistry, surgery, physiology, genetics, psychological_features, psychology, radiology, genetics, dentistry, psychiatry, optics, chemistry

WordNet extended Domains, 117,536 synset-weight pairs
5 most probable for adenoma

14235793-n	0.00010198	medicine
14235793-n	0.00005412	veterinary
14235793-n	0.00003494	anatomy
14235793-n	0.00001745	radiology
14235793-n	0.00001649	cycling

Step to integrate WordNet

Step 4. Map to CSI (Coarse Sense Inventory)

- 45 high-level semantic labels (Coarse Sense Inventory)
- High-level semantics is domain-based
- Manually crafted labels

We select label **HEALTH_AND_MEDICINE_**

C. Lacerra, M. Bevilacqua, T. Pasini, and R. Navigli. 2020. CSI: A coarse sense inventory for 85% word sense disambiguation. In Proceedings of the 34th Conference on Artificial Intelligence, pages 8123–8130. AAAI Press.

Synsets Clinical Domain

Total Wikidata items	26,953
CSI	3,133
WordNet clinical domains	3,396
Total clinical domain only	2,398

Multilingual Central Repository

Step 5. Map to Multilingual WordNet (MCR)

- English
- Spanish
- Catalan
- Basque
- Galician
- Portuguese

```
"WordNet": [  
  {  
    "WordNet 3.1": "14259275-n",  
    "WordNet 3.0": "14235793-n",  
    "CSI": "HEALTH_AND_MEDICINE_",  
    "WordNet Domain": "medicine",  
    "sense": "adenoma%1:26:00::",  
    "MCR synset": [  
      {  
        "en": "a benign epithelial tumor of glandular origin",  
        "es": "tumor epitelial benigno de origen glandular",  
        "pt": "um tumor epitelial benigno de origem glandular",  
        "gl": "",  
        "eu": "",  
        "ca": ""  
      }  
    ]  
  }  
]
```

Result

```
1 {
2   "source_type": "UMLS",
3   "source_id": "C0001430",
4   "status": "OK",
5   "UMLS_CUI": [
6     {
7       "id": "C0001430",
8       "description": "Adenoma"
9     },
10    {
11      "id": "C0001430",
12      "description": "Adenoma, NOS"
13    },
14    {
15      "id": "C0001430",
16      "description": "[M]Adenoma NOS"
17    },
18    {
19      "id": "C0001430",
20      "description": "[M]Adenomas"
21    },
22    {
23      "id": "C0001430",
24      "description": "Benign adenoma"
25    },
26    {
27      "id": "C0001430",
28      "description": "[M]Adenoma NOS (morphologic abnormality)"
29    },
30    {
31      "id": "C0001430",
32      "description": "Adenoma, no subtype (morphologic abnormality)"
33    },
34    {
35      "id": "C0001430",
36      "description": "Adenoma, no subtype"
37    },
38    {
39      "id": "C0001430",
40      "description": "Benign adenomatous neoplasm (disorder)"
41    },
42    {
43      "id": "C0001430",
44      "description": "Benign adenomatous neoplasm"
45    }
46  ],
47  "SNOMED_CT_EN": [
48    {
49      "id": "443416007",
50      "description": "Benign adenomatous neoplasm (disorder) Benign adenomatous neoplasm Adenoma Benign adenoma"
51    },
52    {
53      "id": "32048006",
54      "description": "Adenoma Adenoma, NOS Adenoma, no subtype (morphologic abnormality) Adenoma, no subtype"
55    },
56    {
57      "id": "189579004",
58      "description": "[M]Adenoma NOS [M]Adenoma NOS (morphologic abnormality)"
59    },
60    {
61      "id": "189578007",
62      "description": "[M]Adenomas &/or adenocarcinomas [M]Adenomas and adenocarcinomas [M]Adenomas [M]Adenocarcinomas [M]Adenomas &/or adenocarcinomas (disorder)"
63    }
64  ]
65 }
```

```
64 },
65 "SNOMED_CT_ES": [
66   {
67     "id": "32048006",
68     "description": "adenoma"
69   },
70   {
71     "id": "32048006",
72     "description": "morfología: adenoma, no tipificado (anomalía morfológica)"
73   }
74 ],
75 "ICD10CM_ES": [
76   {
77     "id": "D36.9",
78     "description": "Neoplasia benigna, localización no especificada"
79   }
80 ],
81 "ICD10PCS_ES": [],
82 "MESH": [
83   {
84     "id": "D000236",
85     "description": "Adenoma, Basal Cell"
86   },
87   {
88     "id": "D000236",
89     "description": "Adenoma, Follicular"
90   },
91   {
92     "id": "D000236",
93     "description": "Adenoma, Microcystic"
94   }
95 ],
96 "wikidata_item_url": [
97   "http://www.wikidata.org/entity/Q272741"
98 ],
99 "wikipedia_article_url": [
100   {
101     "arwiki": "https://ar.wikipedia.org/wiki/_>"
102   },
103   {
104     "zhwiki": "https://zh.wikipedia.org/wiki/"
105   }
106 ],
107 "WordNet": [
108   {
109     "WordNet 3.1": "14259275-n",
110     "WordNet 3.0": "14235793-n",
111     "CSI": "HEALTH_AND_MEDICINE_",
112     "WordNet Domain": "medicine",
113     "sense": "adenoma&l:26:00:",
114     "MCR synset": [
115       {
116         "en": "a benign epithelial tumor of glandular origin",
117         "es": "tumor epitelial benigno de origen glandular",
118         "pt": "um tumor epitelial benigno de origem glandular",
119         "gl": "",
120         "eu": "",
121         "ca": ""
122       }
123     ]
124   }
125 ]
126 }
```

Clinical concepts mapped

item	label	MESH	CUI	ICD-10	SNOMED-CT	WN 3.1	WN 3.0	sense	domain	CSI
Q272741	adenoma	D000236	C0001430	D35.0	32048006	14259275-n	14235793-n	adenoma%1:26:00::	medicine	HEALTH_AND_MEDICINE_
Q272741	adenoma	D000236	C0334389	D35.2	32048006	14259275-n	14235793-n	adenoma%1:26:00::	medicine	HEALTH_AND_MEDICINE_
Q7365	muscle organ	D009132	C0026845			05296796-n	05289297-n	musculus%1:08:00::	health	BIOLOGY_
Q84133	myocardium	D009206	C0027061			05398343-n	05391000-n	myocardium%1:08:00::	anatomy	HEALTH_AND_MEDICINE_
Q223102	peritonitis	D010538	C0029823	K65		14376092-n	14352687-n	peritonitis%1:26:00::	medicine	HEALTH_AND_MEDICINE_

Table 5: Examples of WordNets mapped with clinical IDs, WordNet domains and CSI.

Index

1. Introduction
2. Motivation
3. Background
4. Mapping method
5. Corpora Annotation
6. Conclusion

- **CodiEsp** ES: clinical narratives
- **E3C Corpus** ES: clinical narratives
- **CT-EBM-SP** ES: clinical trials
- **Mantra** ES: clinical narratives
- **MedMentions** EN: biomedical papers

Corpora annotated with WordNet synsets

Corpus	Tokens	Annotated CUI	Mapped WN	Unique WN
E3C ES (Magnini et al., 2020)	28,815	2,268	422	107
MedMentions (Mohan and Li, 2019)	1,258,847	540,138	24,754	841
Mantra (Kors et al., 2015)	3,492	1,058	117	62
CT-EBM-SP (Campillos-Llanos, 2019)	141,158	23,264	5,786	431
CodiEsp 2020 (Miranda-Escalada et al., 2020)	401,010	32,902	11,464	399

Table 4: Number of tokens annotated with both WN synsets and clinical IDs using mapping of UMLS CUI to WN synsets.

Corpora annotated with CSI



Obliteration of the upper extraconal fat leads to inferior displacement and ocular proptosis.

La obliteración de la grasa extraconal superior condiciona desplazamiento inferior y proptosis ocular.

Corpora annotated with CSI



Obliteration of the upper extraconal fat leads to inferior displacement and ocular proptosis.

La **obliteración de la grasa** extraconal superior condiciona desplazamiento inferior y **proptosis ocular**.

Corpora annotated with CSI

Obliteration of the upper extraconal fat leads to inferior displacement and ocular proptosis.

La **obliteración de la grasa** extraconal superior condiciona desplazamiento inferior y **proptosis ocular**.

CUI C0015668

SNOMED 238888007 Fat necrosis

MESH D005218 Necrosis, Fat

CUI C0015300

WN 3.1 14313017-n

exophthalmos%1:26:00::

WN 3.0 14336444-n

MESH D005094 exophthalmos

ICD-10 H05: Constant
exophthalmos

optics

HEALTH_AND_MEDICINE_

Index

1. Introduction
2. Motivation
3. Background
4. Mapping method
5. Corpora Annotation
6. Conclusion

Conclusions



- WordNets in various languages are integrated into ClinIDMap
- The mapping tool is scalable for different languages and datasets
- Annotated corpora for sequence labelling models

<https://github.com/Vicomtech/ClinIDMap>

We release

- Source code
- The resulting table of Wikidata to WordNets mapping

**THANK YOU
ESKERRIK ASKO
GRACIAS**



www.vicomtech.org



ezotova@vicomtech.org

REFERENCES

- Campillos-Llanos, L. (2019). First Steps towards Building a Medical Lexicon for Spanish with Linguistic and Semantic Information. pages 152–164, August
- Donnelly, K. et al. (2006). SNOMED-CT: The advanced terminology and coding system for eHealth. *Studies in health technology and informatics*, 121:279.
- Fung, K. W. and Xu, J. (2012). Synergism between the Mapping Projects from SNOMED CT to ICD-10 and ICD-10-CM. *AMIA ... Annual Symposium proceedings. AMIA Symposium*, 2012:218–227.
- Kors, J. A., Clematide, S., Akhondi, S. A., van Mulligen, E. M., and Rebholz-Schuhmann, D. (2015). A multilingual gold-standard corpus for biomedical concept recognition: the Mantra GSC. *Journal of the American Medical Informatics Association*, 22(5):948–956, 05.
- Magnini, B., Altuna, B., Lavelli, A., Speranza, M., and Zanoli, R. (2020). The E3C Project: Collection and Annotation of a Multilingual Corpus of Clinical Cases. (2019). *MedMentions: A Large Biomedical Corpus Annotated with UMLS Concepts*. ArXiv, abs/1902.09476.
- Miranda-Escalada, A., Gonzalez-Agirre, A., Armengol-Estape, J., and Krallinger, M. (2020). Overview of automatic clinical coding: annotations, guidelines, and solutions for non-english clinical cases at Codiesp track of CLEF eHealth 2020. Mohan, S. and Li, D. (2019). *MedMentions: A Large Biomedical Corpus Annotated with UMLS Concepts*. ArXiv, abs/1902.09476.
- Mohan, S. and Li, D. (2019). *MedMentions: A Large Biomedical Corpus Annotated with UMLS Concepts*. ArXiv, abs/1902.09476.
- Elena Zotova, Montse Cuadros, and German Rigau. 2022. ClinIDMap: Towards a Clinical IDs Mapping for Data Interoperability. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 3661–3669, Marseille, France. European Language Resources Association.
- Barry Smith and Christiane Fellbaum. 2004. Medical wordnet: A new methodology for the construction and validation of information resources for consumer health. In *Proceedings of the 20th International Conference on Computational Linguistics, COLING '04*, page 371es, USA. Association for Computational Linguistics.
- C. Lacerra, M. Bevilacqua, T. Pasini, and R. Navigli. 2020. *CSI: A coarse sense inventory for 85% word sense disambiguation*. In *Proceedings of the 34th Conference on Artificial Intelligence*, pages 8123–8130. AAAI Press.