

BBMRI.at

Biobanking and
BioMolecular resources
Research Infrastructure
Austria

Optimizing Biobank Data Retrieval by Using IT-Based Lexical and Semantic Methods

Philipp Hofer-Picout, Sabrina Neururer, Georg Göbel

philipp.hofer@i-med.ac.at



Funded by GZ 10.470/0016-II/3/2013

Table of Contents



BBMRI.at

Biobanking and
BioMolecular resources
Research Infrastructure
Austria

Background

Objectives

Methods

Results

Query Challenges

Lexical & Semantic Methods

Prototype Implementation

Proof-of-Concept

Discussion & Outlook

Query Heterogeneity Problem

No Results Found.

No Results Found.

wilms tumor SEARCH

kidney cancer SEARCH

No Results Found.

ICD10:C64 SEARCH



BIMS*



(i) Sample set
Material type: FFPE, Frozen Tissue
Number of samples: 2
Diagnosis: **nephroblastoma**

...

*Biobank Information Management System

Facilitate Data Accessibility

→ Optimize retrieval of desired information



(i) Sample set
Material type: FFPE, Frozen Tissue
Number of samples: 2
Diagnosis: **nephroblastoma**



(i) Sample set
Material type: FFPE, Frozen Tissue
Number of samples: 2
Diagnosis: **nephroblastoma**



(i) Sample set
Material type: FFPE, Frozen Tissue
Number of samples: 2
Diagnosis: **nephroblastoma**



BIMS

Main Steps

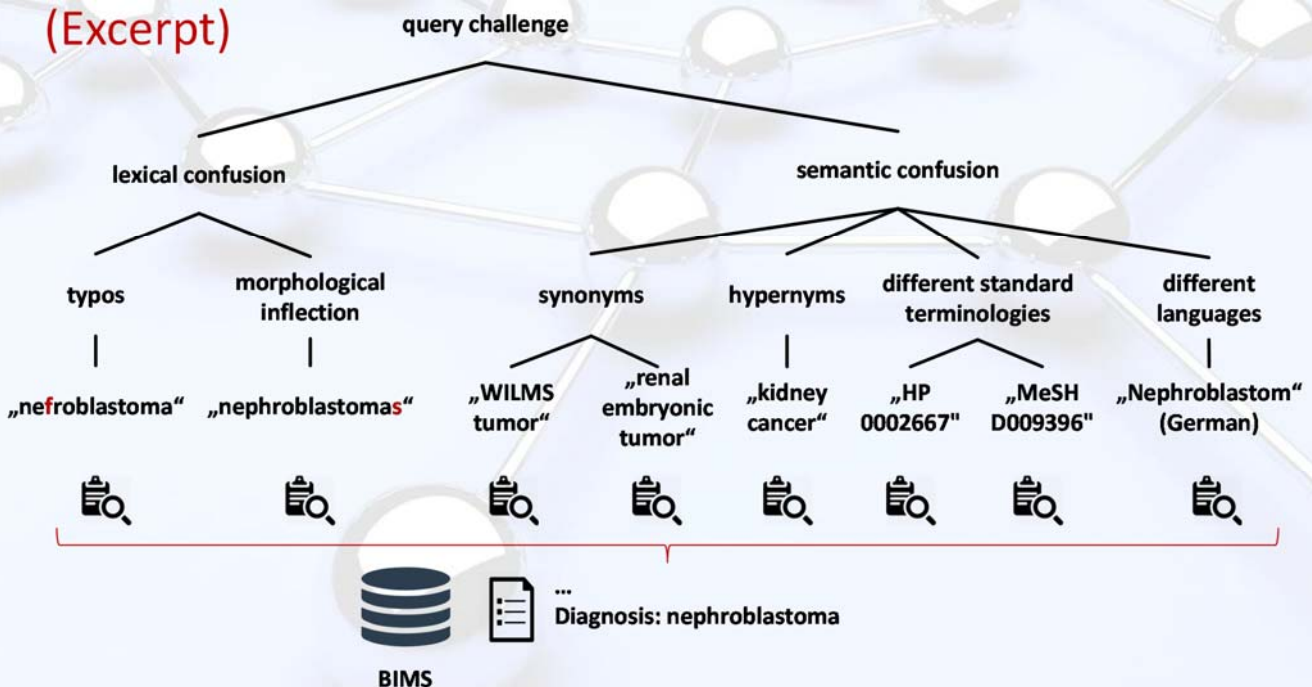
- I. Identify & classify different query challenges^{1,2}
- ↓
- II. Find appropriate methods for all query scenarios
- ↓
- III. Prototype implementation
- ↓
- IV. Proof-of-concept

[1] Fidel, R. (1985). Moves in online searching. Online Review. <https://doi.org/10.1108/eb024176>

[2] Liu, C., Gwizdka, J., Liu, J., Xu, T., & Belkin, N. J. (2010). Analysis and evaluation of query reformulations in different task types. In Proceedings of the ASIST Annual Meeting. <https://doi.org/10.1002/meet.14504701214>

I. Query Challenges

(Excerpt)



II. Lexical & Semantic Methods

- **Lexical confusion**³
 - Vector Space Model
 - Phonetic Analysis/ Ngram/ Stemming (english)
 - Approximate String Matching (Levenstein distance)
- **Semantic confusion**
 - RDF Data Representation/ Graph Database⁴
 - (Bio-)Medical Standard Terminologies⁵

[3] M.S. Divya, S.K. Goyal, ElasticSearch: An advanced and quick search technique to handle voluminous data, International Journal of Advanced Computer Technology, 6 (2013), 171-175.

[4] <https://www.blazegraph.com>, last accessed on April 2018.

[5] Tirmizi SH, Aitken S, Moreira DA, Mungall C, Sequeda J, Shah NH, et al. Mapping between the OBO and OWL ontology languages. J Biomed Semantics. BioMed Central Ltd; 2011;2 (Suppl 1):S3.

III. Prototype Implementation

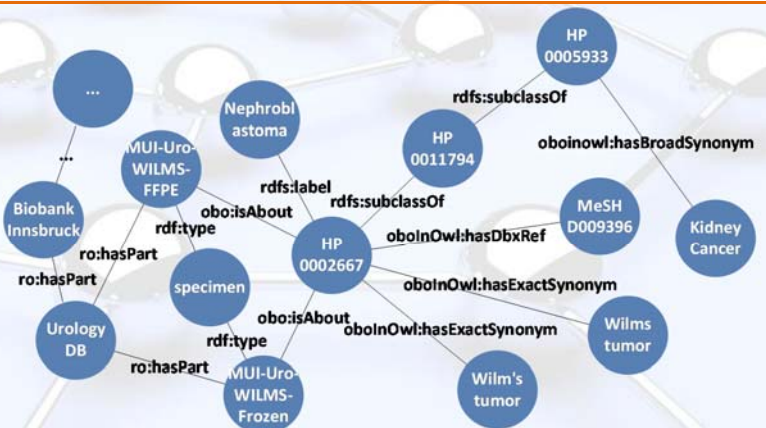


(I) Sample set⁶
 Material type: FFPE, Frozen Tissue
 Number of samples: 2
 Diagnosis: nephroblastoma

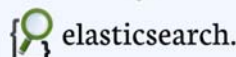


Disease ontologies:
 DOID, Orphanet, ICD10-CM,
 HPO, MESH

**CREATE
 RDF
 DATA**



**CREATE
 INDEX**



```
select ?uri (STR(?l) AS ?label)
WHERE {
  ?uri a owl:Class .
  { ?uri rdfs:label ?l }
  UNION { ?uri oboinowl:hasExactSynonym ?l }
  UNION { ?uri oboinowl:hasRelatedSynonym ?l }
  UNION { ?uri oboinowl:hasBroadSynonym ?l }
  UNION { ?uri skos:prefLabel ?l }
  UNION { ?uri skos:notation ?l } .
}
```

[6] Hofer P, Fiegl H, Angerer J, Mueller- E. A Concept of a MIABIS based Register of Biosample Collections at the Medical University of Innsbruck. 2014;

IV. Proof-Of-Concept

Query class	Input String	Annotation used	Record found†	Applied method
misspelling	nefroblastom	HPO ^a 0002667	Yes	Approximate string matching
Hypernym	Cancer of the kidneys	HPO 0002667	Yes	hierarchical relationship
Synonym	Wilms tumor	HPO 0002667	Yes	synonym relationship
coding standard	MeSH ^b D009396	HPO 0002667	Yes	cross reference HPO → MeSH

†Expected record:



(i) Sample set
Material type: FFPE, Frozen Tissue
Number of samples: 2
Diagnosis: **nephroblastoma**

22.09.2018

^aHuman Phenotype Ontology

^bMedical Subject Headings

9

Discussion & Outlook

- Both stemming and phonetic algorithms are more robust against typos and morphological variations
- Cross-references allow searching across different terminologies
- (Bio-)Medical standard terminologies must be valid regarding knowledge representation
- Combining different techniques would likely be the most rewarding
- A long-term evaluation about frequency of different query scenarios in productive environments will be necessary

22.09.2018

10

Literature

- [1] Fidel, R. (1985). Moves in online searching. Online Review. <https://doi.org/10.1108/eb024176>
- [2] Liu, C., Gwizdka, J., Liu, J., Xu, T., & Belkin, N. J. (2010). Analysis and evaluation of query reformulations in different task types. In Proceedings of the ASIST Annual Meeting. <https://doi.org/10.1002/meet.14504701214>
- [3] M.S. Divya, S.K. Goyal, ElasticSearch: An advanced and quick search technique to handle voluminous data, International Journal of Advanced Computer Technology, 6 (2013), 171-175.
- [4] <https://www.blazegraph.com>, last accessed on April 2018.
- [5] Tirmizi SH, Aitken S, Moreira DA, Mungall C, Sequeda J, Shah NH, et al. Mapping between the OBO and OWL ontology languages. J Biomed Semantics. BioMed Central Ltd; 2011;2 (Suppl 1):S3.
- [6] Hofer P, Fiegl H, Angerer J, Mueller- E. A Concept of a MIABIS based Register of Biosample Collections at the Medical University of Innsbruck. 2014;

Thank you for your attention !



Georg Goebel

Principal Investigator
Local biobank coordinator
Biobank-IT contact
Medical Statistics, Informatics and Health Economics
georg.goebel@i-med.ac.at



Sabrina Neururer

Local project coordination
Medical Statistics, Informatics and Health Economics
sabrina.neururer@i-med.ac.at



Philipp Hofer-Picout

Researcher
Medical Statistics, Informatics and Health Economics
philipp.hofer@i-med.ac.at

Contact

biobank@i-med.ac.at
Tel. +43 512 9003 70912
Fax +43 512 9003 73922

**IBOBANK
INNSBRUCK**

**Σ MEDICAL
STATISTICS
INFORMATICS
HEALTH ECONOMICS**

BBMRI.at
Neue Stiftingtalstraße 2/B/6
A-8010 Graz
Austria
www.bbmri.at