

# Medical Concept Representation - the Years Beyond 2000

Laszlo Balkanyi<sup>a</sup>, Stefan Schulz<sup>b,c</sup>, Ronald Cornet<sup>d,e</sup>, Olivier Bodenreider<sup>f</sup>

<sup>a</sup>European Centre for Disease Prevention and Control, Stockholm, Sweden, <sup>b</sup>Medical University of Graz, Austria, <sup>c</sup>Freiburg University Medical Center, Freiburg, Germany, <sup>d</sup>Academic Medical Center, Amsterdam, The Netherlands, <sup>e</sup>Linköping University, Linköping, Sweden, <sup>f</sup>National Library of Medicine, Bethesda, USA

## Abstract

- understanding the state of the art in the context of "medical concept representation"
- a descriptive study based on bibliometrics, simple text mining and a social media survey
- results support the general understanding that the focus of research has moved toward medical ontologies
- socially active researchers mention the OBO Foundry, SNOMED, and UMLS as key resources
- text mining of most cited literature identifies single noun phrases as "health", "information", "clinical", "knowledge", "ontology", "case", "data", "semantic(s)", "concept" and "representation" as leading denominators of the field
- terms as "ontology" and "semantic(s)" have gained more significance in the last decade
- there is a paradigm shift according to both the socially active group of researchers and bibliometric data, comparing citation ranks of the nineties and the recent decade support this opinion

## Introduction

The goal of the exploratory study: to understand the state of the art in the broad contextual research area of "medical concept representation" originating in the 1990s

- Influencing factors:**
- advances in medical information science,
  - terminologies, ontology development,
  - accessibility of networked computing

**Effect:**

- significant growth of research, development practical implementation in this area [1, 2].

**Background:**  
The study was initiated by the IMIA Medical Concept Representation Working Group (MCR WG).

The WG was one of the most influential bodies in the late eighties and the nineties, publishing regularly overviews of the domain, last in-depth, analytic study dated to 2006 [3].

- Current study based on:**
- bibliographic measures,
  - simple on-line text mining tools,
  - a social media survey

## Materials and Methods

### Selecting sources of information on most influential papers :

- Scopus term analyzer [6] for a time line for catch phrase "medical concept representation" used in titles of publications
- Ultimate Research Assistant [7] to extract contextual environment
- Wordle tag cloud tool to visualize the context [8]
- Identification of authors of **ten most influential papers** by using seven sources: 1) Web of Science, 2) Scopus, 3) Embase, 4) PubMed, 5) Google Scholar, 6) Cochrane Library, 7) British Library on-line catalogue.
- Boolean search expression "concept representation" AND (medical OR medicine) AND (knowledge OR information) to filter non-relevant results
- two periods, 1988 – 1999 and 2000 – 2012, were compared by text-mining titles of publications

### Setting up a social network survey:

- Primary source: LinkedIn group of the IMIA Medical Concept Representation WG, which has over 50 members with widely various backgrounds.
- Survey open from August 2012 to the end of October 2012.
- Secondary sources were additional LinkedIn Groups in broader domain where the survey was also published.
- Lists of most relevant papers shared by Datagly [9] and a tailored Google Docs document.

### Reconciling bibliometrics and the survey data:

- noun phrase frequencies compared using text mining results of paper titles
- text mining performed with Textalyser [10]

## Results 1

Figure 1- time line analytics results for the phrase "medical concept representation": the use of the catch phrase in the nineties and in the first decade of the new century. The timeline reveals that studies using this title phrase were done in mostly the nineties.

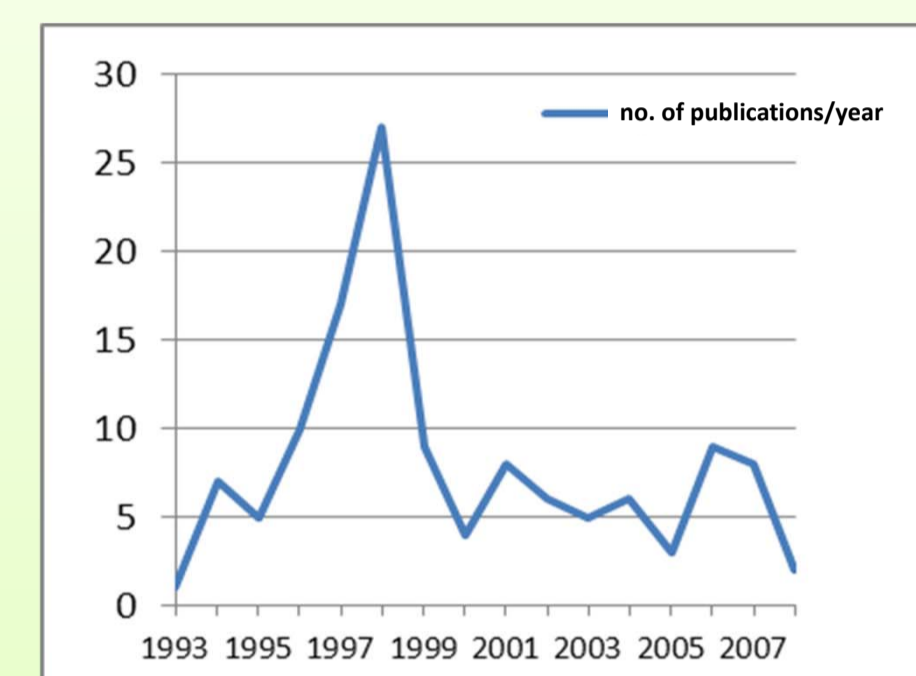


Figure 2- Tag cloud generated from result of catch phrase search reveals the wide contextual domain related to the catch phrase "medical concept representation".



## Results 2

Table 1 – Set of thirty most influential authors of the period 1988-1999

Ranks were calculated by adding positional scoring of papers in the various lists, weighted in favour of multiple appearances of authors.

score	authors (1-15)	score	authors (16-30)
83.2	Clintino JJ	22.8	Rosse C,
43.2	Oliver DE	21	Miller R.A.
39.2	Burd RH	21	Resnik AM
39.2	Scherrer JR	19.2	Mason MA
35	Rector AL	19.2	Nowlan WA
33.6	Bell DS	18.2	Wagner JC
33.6	Shahar Y	18	Bailey RR
33.6	Shoaff EH	18	Bauer BA
30.8	Fieschi M	18	Elkin PL
30.8	Huff SM	16.8	Cham CG
30.8	Joubert M	15.6	Schoorman HM
30.8	Valde F	15.6	Barnett GO
26.6	Johnson SB	15.6	Horrocks I
22.8	Evans DA	15.6	Humphrey BL
22.8	Hersh WR	15.6	Lindberg DAB

Table 2 – Set of most cited authors, period 2000 – 2012, in the composite contextual domain of authors on medical concept representation. (The full citation data was compared, not limited to citing the exact search phrase. )

By comparing author's citation data over the composite contextual domain a broader coverage and less bias toward particular phrasing is achieved. Only about ten percent of mostly cited authors of the nineties remained in the list of the recent period.

authors (1-15)	cited	authors (16-30)	cited
Smith B	12529	Noy NF	21195
Roberts A	42871	Nathan SE	20866
Stevens R	62715	Jaffe H	19204
Horrocks I	60177	Wise C	17530
Van Harmelen F	58626	Lustier Y	14802
Fensel D	58401	Coronado S	14105
Zakhe LA	49420	Srinivasan C	6976
Goble C	46984	Srinivas N	6584
Hollman KM	42858	Yao YY	3516
Dickler S	41092	Shapiro L	3407
Friedman C	40473	Harris FW	3211
Fal SK	31200	Mejino JR	2925
Mason MA	28181	Haber MW	2525
Aspden P	23482	Shiu SK	1611
Rosse C	22553	Steinman F	1074

## Results 3

Table 3 - A word frequency analysis of the titles of all 'medical concept representation' papers of the period 1988-1999 shows the ten most frequently used single noun phrases and meaningful two-word phrases.

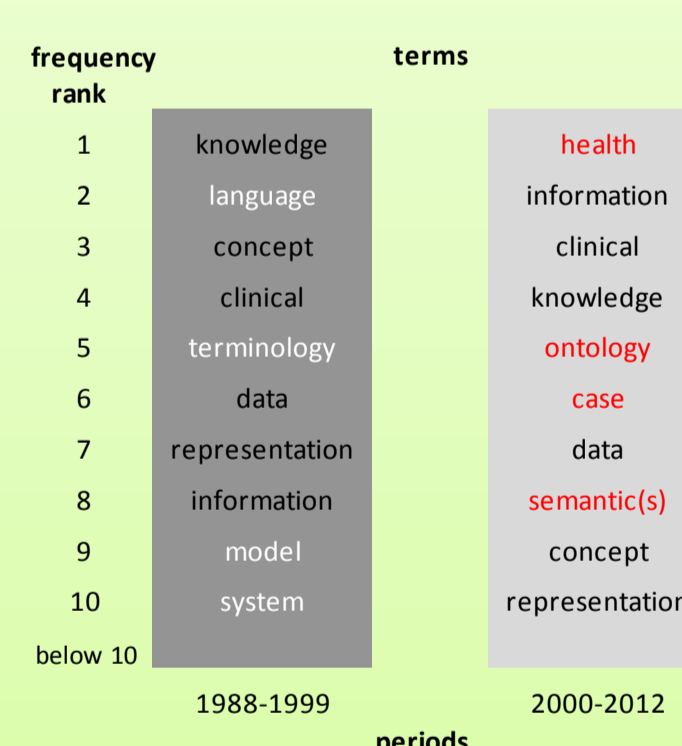
rank	single terms	two-word phrases
1	knowledge	medical language
2	language	natural language
3	concept	case based
4	clinical	knowledge representation
5	terminology	knowledge acquisition
6	data	language processing
7	representation	medical concept
8	information	medical terminology
9	model	structured data
10	system	concept representation

Table 4 – List of most frequently used single noun phrases and two word phrases of the period 2000-2012 in the domain of medical concept representation. All red entries are new terms or phrases in the titles of papers published, compared to the list of phrases of the 1988-1999 period.

rank	single terms	two-word phrases
1	health	health informatics
2	information	electronic health
3	clinical	natural language
4	knowledge	concept based
5	ontology	decision support
6	case	language processing
7	data	concept representation
8	semantic(s)	medical language
9	concept	medical informatics
10	representation	description logic

## Results 4

Table 5 - title term changes : outgoing terms are in white. New title terms are red. The contextual domain of formal medical concept representation was broadened. Prominence of terms "semantics" and "ontology" shows a new paradigm. Terms "language", "model" and "terminology" disappeared - these more differentiated areas branched off from the previously common roots.



## Results 5

Table 6 - Results of the survey of socially active researchers got over 40 responses. The idea of using ontologies has clearly become the central paradigm. The list of top 20 titles of multiple occurrences are presented. Major research areas and resources in these lists include the OBO Foundry, the Gene Ontology, various versions of Systematized Nomenclature of Medicine (SNOMED) and the Unified Medical Language System (UMLS).

Rank	Title
1	The OBO Foundry: coordinated evolution of ontologies to support biomedical data integration (649) [13]
2	The Unified Medical Language System (UMLS): integrating biomedical terminology (709) [11]
2	A reference ontology for biomedical informatics: the Foundational Model of Anatomy (76) [14]
2	Relations in biomedical ontologies (62) [15]
2	Disorders for controlled medical vocabularies in the twenty-first century (457) [16]
2	Clinical terminology: why is it so hard? (242) [17]
2	From concepts to clinical reality: an essay on the benchmarking of biomedical terminologies (37) [18]
2	Web-Center: detecting public health rumors with a Web-based text mining system (6) [19]
2	An ontology of epidemiological terms (ref website) [20]
2	Gene ontology: tool for the unification of biology (1009) [21]
2	Switzerland Ontologies with DOLCE (66) [22]
2	The medical dictionary for regulatory activities (MDRA) (20) [23]
2	SNOMED clinical terms: overview of the development process and project status (150) [24]
2	Towards a reference terminology for ontology research and development in the biomedical domain (92) [25]
2	Methods in biomedical ontology (2) [26]
2	Ontology-based error detection in SNOMED-CT (82) [27]
2	Fuzzy Health, Illness, and Disease (66) [28]
2	Medicina bioinformatica: experimental processes with (MIS) (29)
2	Integrating epidemiology into the Semantic Web (1) [30]
2	A Dictionary of Epidemiology (book) [31]

## Discussion and conclusions

### Methodology issues:

- this is an ad hoc methodology not aiming at a new scientometric index or a generalized, reusable tool to measure or assess other areas
- combining on-line library databases, bibliometric services and text mining tools resulted in study-focused tool sets while large size of screened sources alleviate possible bias
- using the notion "concept" decreased probably due to the following factors: (A) propagation of the paradigm of ontological realism and critique of "conceptualism"[33], (B) substitution of the notion "concept" by the notion "class" in the Semantic Web and description logics community [34], (C) by the obvious ambiguity of the word itself [35]

### Conclusions:

- the new millennium has coincided with a change in the focus of research
- the new paradigm of Semantic Web and results of ontology theory and practice have become new anchors [36].
- the central role of the term "concept" has been gradually abandoned
- the new paradigm could be illustrated with the following exemplars:
  - capture of medical information and knowledge leverages (standard based) ontologies
  - open reference resources for content are developed collaboratively and are increasingly shared, reused
  - web enabled standards help achieve transparent results.

### Acknowledgments

This research was supported in part by the Intramural Research Program of the National Institutes of Health (NIH), National Library of Medicine (NLM).

Reference list: <http://goo.gl/EJZrg> Mail to: [laszlo.balkanyi@ecdc.europa.eu](mailto:laszlo.balkanyi@ecdc.europa.eu)