



Contents lists available at ScienceDirect

# Journal of Biomedical Informatics

journal homepage: [www.elsevier.com/locate/yjbin](http://www.elsevier.com/locate/yjbin)



## Context-driven automatic subgraph creation for literature-based discovery

Delroy Cameron<sup>a,\*</sup>, Ramakanth Kavuluru<sup>b</sup>, Thomas C. Rindflesch<sup>c</sup>, Amit P. Sheth<sup>a</sup>,  
Krishnaprasad Thirunarayan<sup>a</sup>, Olivier Bodenreider<sup>c</sup>

<sup>a</sup> Ohio Center of Excellence in Knowledge-Enabled Computing (Kno.e.sis), Wright State University, Dayton, OH 45435, USA

<sup>b</sup> Division of Biomedical Informatics, University of Kentucky, Lexington, KY 40506, USA

<sup>c</sup> National Library of Medicine, 8600 Rockville Pike, Bethesda, MD 20894, USA

### ARTICLE INFO

**Article history:**  
Received 16 September 2014  
Accepted 25 January 2015  
Available online xxxxx

**Keywords:**  
Literature-based discovery (LBD)  
Graph mining  
Path clustering  
Hierarchical agglomerative clustering  
Semantic Similarity  
Semantic relatedness  
Medical Subject Headings (MeSH)

### ABSTRACT

**Background:** Literature-based discovery (LBD) is characterized by uncovering hidden associations in non-interacting scientific literature. Prior approaches to LBD include use of: (1) domain expertise and structured background knowledge to manually filter and explore the literature, (2) distributional statistics and graph-theoretic measures to rank interesting connections, and (3) heuristics to help eliminate spurious connections. However, manual approaches to LBD are not scalable and purely distributional approaches may not be sufficient to obtain insights into the meaning of poorly understood associations. While several graph-based approaches have the potential to elucidate associations, their effectiveness has not been fully demonstrated. A considerable degree of *a priori* knowledge, heuristics, and manual filtering is still required.

**Objectives:** In this paper we implement and evaluate a context-driven, automatic subgraph creation method that captures multifaceted complex associations between biomedical concepts to facilitate LBD. Given a pair of concepts, our method automatically generates a ranked list of subgraphs, which provide informative and potentially unknown associations between such concepts.

**Methods:** To generate subgraphs, the set of all MEDLINE articles that contain either of the two specified concepts (A, C) are first collected. Then binary relationships or assertions, which are automatically extracted from the MEDLINE articles, called *semantic predications*, are used to create a labeled directed *predications graph*. In this predications graph, a *path* is represented as a sequence of semantic predications. The hierarchical agglomerative clustering (HAC) algorithm is then applied to cluster paths that are bounded by the two concepts (A, C). HAC relies on implicit semantics captured through Medical Subject Heading (MeSH) descriptors, and explicit semantics from the MeSH hierarchy, for clustering. Paths that exceed a threshold of semantic relatedness are clustered into subgraphs based on their *shared context*. Finally, the automatically generated clusters are provided as a ranked list of subgraphs.

**Results:** The subgraphs generated using this approach facilitated the rediscovery of 8 out of 9 existing scientific discoveries. In particular, they directly (or indirectly) led to the recovery of several *intermediates* (or B-concepts) between A- and C-terms, while also providing insights into the meaning of the associations. Such meaning is derived from predicates between the concepts, as well as the provenance of the semantic predications in MEDLINE. Additionally, by generating subgraphs on different thematic dimensions (such as *Cellular Activity*, *Pharmaceutical Treatment* and *Tissue Function*), the approach may enable a broader understanding of the nature of complex associations between concepts. Finally, in a statistical evaluation to determine the *interestingness* of the subgraphs, it was observed that an arbitrary association is mentioned in only approximately 4 articles in MEDLINE on average.

**Conclusion:** These results suggest that leveraging the implicit and explicit semantics provided by manually assigned MeSH descriptors is an effective representation for capturing the underlying *context* of complex associations, along multiple thematic dimensions in LBD situations.

© 2015 Published by Elsevier Inc.

\* Corresponding author. Tel.: +1 937 775 5213; fax: +1 937 775 5133.  
E-mail address: [delroy@knoesis.org](mailto:delroy@knoesis.org) (D. Cameron).

70 **1. Introduction**

71 Literature-based discovery (LBD) refers to the process of uncov-  
72 ering hidden connections that are implicit in scientific literature.  
73 Numerous hypotheses have been generated from scientific litera-  
74 ture, using the LBD paradigm, which influenced innovations in  
75 diagnosis, treatment, preventions, and overall public health. The  
76 notion of LBD was proposed by Don R. Swanson (1924–2012) in  
77 1986, through the well-known *Raynaud Syndrome–Dietary Fish Oils*  
78 *Hypothesis (RS-DFO)* [1]. By reading the titles of more than 4000  
79 MEDLINE articles, Swanson serendipitously discovered that *Dietary*  
80 *Fish Oils (DFO)* lower *Blood Viscosity*, reduce *Platelet Aggregation*  
81 and inhibit *Vascular Reactivity* (specifically *Vasoconstriction*). Concomi-  
82 tantly, he observed that a reduction in both *Blood Viscosity* and  
83 *Platelet Aggregation*, as well as the inhibition of *Vasoconstriction*,  
84 appeared to prevent *Raynaud Disease*; a circulatory disorder that  
85 causes periods of severely restricted blood flow to the fingers  
86 and toes [2]. Swanson therefore postulated that “*dietary fish oil*  
87 *might ameliorate or prevent Raynaud’s syndrome.*” This hypothesis  
88 was clinically confirmed by DiGiacomo et al. [3] in 1989.

89 Swanson’s discovery is interesting because explicit associations  
90 between *DFO* and these intermediate concepts (i.e., *Blood Viscosity*,  
91 *Platelet Aggregation* and *Vasoconstriction*) had long existed in the  
92 literature [4–8]. Likewise, explicit associations between the inter-  
93 mediates and *RS* had been well documented [9,2]. The serendipity  
94 in Swanson’s Hypothesis lies in the fact that no explicit  
95 associations linking *DFO* and *RS* directly had been previously  
96 articulated in a single document.

97 To develop this hypothesis, Swanson performed a Dialog® Sci-  
98 search using *Raynaud* and *Fish Oil* terms, on titles and abstracts  
99 of MEDLINE and Embase (Excerpta Medica) citations, in November  
100 1985. There were approximately 1000 articles in the *Raynaud* set  
101 and 3000 in the *Fish Oil* set. He found that only four articles among  
102 a reduced set of 489 articles (after filtering), contained cross-refer-  
103 ences spanning both sets. Among these four articles, only two arti-  
104 cles [10,11] discussed relevant aspects of *RS* with *DFO*; although  
105 not in the context of Swanson’s discovery. Swanson speculated  
106 that this phenomenon of logically related but noninteracting litera-  
107 tures alludes to the existence of *undiscovered public knowledge* [1].  
108 Logically related information fragments may exist in the literature,  
109 but may have never been connected, or fully elucidated. He subse-  
110 quently exploited his awareness of the existence of such undiscovered  
111 associations and investigated several other scenarios (three  
112 with Smalheiser [12–14]) that later led to new scientific discover-  
113 ies [15,16]. Swanson grounded his observations in a paradigm now  
114 commonly known as the *ABC model* [1] for LBD, which is an integral  
115 part of LBD research, facilitating the generation of several hypoth-  
116 eses [1,15,16,12–14,17–25].

117 In current biomedical research, while finding unknown inter-  
118 mediates is an important task, domain scientists are often inter-  
119 ested in developing a deeper understanding of causal  
120 relationships and mechanisms of interaction among concepts. For  
121 example, consider the complex scenario depicted in Fig. 1, in which  
122 *Dietary Fish Oils* produce several *Prostaglandins*, including *Prosta-*  
123 *glandin I3 (PGI<sub>3</sub>)* and *Epoprostenol (PGI<sub>2</sub>, the synthetic form of*  
124 *Prostacyclin)*. The latter of these *Prostaglandins (Epoprostenol)* was  
125 known to treat *Raynaud Syndrome*. It was also known to disrupt  
126 *Platelet Aggregation*. Since *Platelet Aggregation* is deemed a cause  
127 of *Raynaud Syndrome*, one can reasonably conclude that a plausible  
128 mechanism by which *Dietary Fish Oils* treat *Raynaud Syndrome* is  
129 through the production of *Prostaglandins*, which actively disrupt  
130 *Platelet Aggregation*.

131 Aside from detecting such causal associations, it is known that  
132 complex associations may exist between concepts, in many  
133 different ways. For example, Fig. 2 shows that *Dietary Fish Oils*

and *Raynaud Syndrome* are associated in at least the following  
three ways: (1) in terms of *Cellular Activity* involving *Blood plate-*  
*lets/Prostaglandins*, as shown in Fig. 2a, (2) through *Pharmaceuticals*  
that contain calcium channel blockers, such as *Nifedipine* and  
*Verapamil*, as shown in Fig. 2b, and (3) through *Lipids/Fatty Acids*  
from *Efamol* and *Evening primrose oil*, as shown in Fig. 2c.

In this paper, we build on our previous approach [26], in which  
we rediscovered and decomposed the *Raynaud Syndrome – Dietary*  
*Fish Oils* discovery. In our previous work, we manually created the  
multi-faceted subgraphs, by grouping together paths of *semantic*  
*predications*. Recall that a semantic predication is a binary relation  
between two concepts, expressed in the form (subject, predicate,  
object). Here, we present a method that uses rich representations  
to automatically create such subgraphs, by leveraging implicit  
and explicit semantics provided by MeSH descriptors.<sup>1</sup> To create  
the subgraphs, we first specify the context of a semantic predication  
and then use it to infer the context of a path. Paths are then clustered  
into coherent subgraphs on multiple thematic dimensions, based on  
their shared context.

The approach requires only three items from the user as input:  
(1) a list of concept labels for source (A) and target (C), (2) the max-  
imum path length *k* of paths to be generated (default *k* = 2, for *ABC*  
associations), and (3) a cut-off date *dt* for articles to be included  
from the scientific literature. If no cut-off date is provided then  
all MEDLINE articles are used. The output of the approach is a  
ranked list of subgraphs *S* – i.e., create a function  $\mathcal{F}: q \rightarrow S$ , where  
 $q = \{A, C, dt, k\}$ .

To facilitate understanding the meaning of associations present  
in the subgraphs, the predicates of the semantic predications and  
their provenance in MEDLINE are provided (see Section 4). Rela-  
tionships that are not explicit in the subgraphs, but are inferred,  
can be explored by composing MEDLINE queries (as we will show).  
The collective use of predicates, provenance and MEDLINE queries  
for knowledge exploration constitute the notion of *discovery*  
*browsing*, introduced by Wilkowski et al. [27] and extended by  
Cairelli et al. [28]. Discovery browsing is enabled when a system  
guides the user through their exploration of the literature in a pro-  
cess of cooperative reciprocity. The “*user iteratively focuses system*  
*output, thus controlling the large number of relationships often gener-*  
*ated in literature-based discovery systems.*”

To assess the efficacy of our approach, two forms of evaluation  
were conducted: (1) an evidence-based evaluation and (2) a statis-  
tical evaluation. The evidence-based evaluation showed that the  
generated subgraphs could facilitate the rediscovery of 8 out of 9  
existing discoveries [1,15,16,12–14,29,30] (not recovered [28]).  
The statistical evaluation showed that an arbitrary association  
occurs only in approximately 4 articles in MEDLINE on average.  
This evaluation determines the *interestingness* of the subgraphs in  
general, as a way to assess whether a domain scientist might be  
interested in an arbitrary subgraph in the first place (see in Section  
4.2). These results suggest that the subgraphs created using our  
approach provide an effective way of finding and elucidating  
poorly understood associations and may be of interest to domain  
scientists. In this paper we make the following specific  
contributions:

1. We develop a novel context-driven subgraph creation method  
for closed LBD (both A and C are known), capable of finding  
complex associations. Our approach is distinct from previous  
approaches, which are mainly based on statistical frequency,  
graph metrics, and specificity.

<sup>1</sup> MeSH is a controlled vocabulary (or thesaurus) of biomedical terms, organized in  
a hierarchical structure – <https://www.nlm.nih.gov/mesh/>.

Download English Version:

<https://daneshyari.com/en/article/6928209>

Download Persian Version:

<https://daneshyari.com/article/6928209>

[Daneshyari.com](https://daneshyari.com)