

# Demonstration of a Literature Based Discovery System based on Ontologies, Semantic Filters and Word Embeddings for the Raynaud Disease-Fish Oil Rediscovery

Toby Reed<sup>1,2</sup>, Vassilis Cutsuridis<sup>2</sup>

<sup>1</sup>Streets Heaver Healthcare Computing, Lincoln, UK

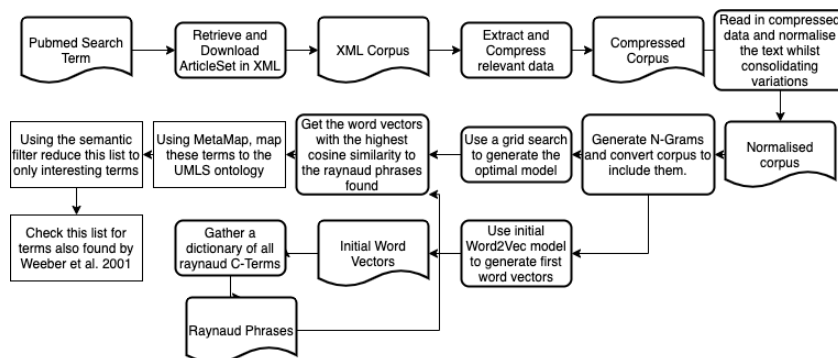
<sup>2</sup>School of Computer Science, University of Lincoln, Lincoln, UK

## Introduction

In the current COVID-19 era there is widespread demand from the pharmaceutical and healthcare industries for more work to be done in the field of reusing compounds in diseases as a method to escape some of the most expensive and time-consuming processes in drug discovery. After the most famous Sildenafil (Viagra) being repurposed from cardiovascular disease to erectile dysfunction, the use of drug repositioning has been shown to have the potential to be beneficial not only to the healthcare facilities and pharmaceutical companies, but also to the everyday consumer particularly if the process to finding and developing cures becomes cheaper, then the actual to consumer cost of treatment will likely decrease [1]. One method for drug repositioning is through Literature Based Discovery (LBD), a powerful text mining approach that harnesses already available scientific knowledge to build bridges between seemingly unrelated islands of knowledge, such as the association of an existing drug to a novel medical condition [1]. LBD is classified into two types: *open* and *closed* discovery. In closed discovery (also known as *hypothesis testing*), the user specifies a pair of topics (A and C) and the objective is to find any unknown, but meaningful connections (the *intermediate* (B) terms) between them [1]. In open discovery (also known as *hypothesis generation*), the user specifies a topic of interest (C) (e.g. a disease or a drug) and the system finds a set of *intermediate* (B) terms directly related to the starting topic of interest [1]. For each of these intermediate terms, the system reiterates the same mechanism to generate a set of *final* (A) terms [1].

## Materials and Methods

A novel LBD system based on Word Embeddings, Statistics, Semantic Filters, and UMLS ontologies was developed to rediscover the Raynaud disease-Fish Oil connection [2] by mining Pubmed titles and abstracts:



Our system's pipeline and corpora mined to discover the Raynaud disease – Fish oil connection [2] can briefly be described as follows:

1. **Corpora:** Different size and specificity corpora of Pubmed titles/abstracts were retrieved for each discovery type (open vs closed). *Open discovery corpora:* (i) Vascular disease, (2) Peripheral vascular disease (PVD), and (3) Raynaud disease. *Closed discovery corpora:* (i) Hematology, (ii) PVD, and (3) Raynaud disease. Raynaud disease is a specific type of vascular disease, but not a type of PVD, which is a sub-type of vascular disease. Raynaud disease is a sub-type of vascular disease, which involves blood, but not per se a sub-type of a hematological disease and neither is PVD.
2. **Pre-processing:** Each retrieved title/abstract of a scientific article was normalized to remove word variations due to capitalization. Any words with the less than three characters was also removed from further processing. All

remaining words were then passed through a Natural Language Toolkit parser to generate bigrams/trigrams of each unigram based on a minimum occurrence count value.

3. A Skip-Gram Word2Vec model [3] was employed with some initial parameter values to generate word vectors for all words and phrases in each corpus.
4. We scanned through all generated word vectors to discover variations of the “raynaud” C-concept (e.g. Raynaud’s disease, Raynaud syndrome, primary Raynaud, etc).
5. We utilised a grid search on the architecture, dimensionality, epoch, learning rate, down-sampling, context window and minimum word count parameters to find the model with the optimum performance in each corpus used.
6. Using the optimally derived Word2Vec model, we repeated STEP 4 to estimate cosine similarity of all B- or A-terms in the corpus with Raynaud variation terms from STEP 3.
7. Placed the most semantically similar terms with the closest cosine similarity, from STEP 5 into a list.
8. Mapped every term from the list via MetaMap [4] to UMLS [5] ontologies. Using a semantic filter we excluded from further analysis all mapped terms which were not semantically related to the semantic types in the filter.
9. These results were then compared to previously found terms to see if our system provided acceptable results.

### Results

We provide below results from both discovery modes (open vs closed). In ‘green’ are the correctly rediscovered concepts, in ‘yellow’ the semantically similar discovered concepts, and in ‘red’ the concepts our system ought to have discovered, but failed to do so:

Open discovery				Closed discovery		
	Vascular	Peripheral Vascular	Raynaud	Raynaud	Peripheral Vascular	Hematology
Blood Viscosity	Green	Green	Yellow	Green	Yellow	Yellow
Platelet Aggregation	Green	Green	Yellow	Green	Yellow	Green
Vascular Reactivity	Green	Yellow	Yellow	Yellow	Yellow	Yellow
Erythrocyte Deformability	Green	Yellow	Yellow	Green	Yellow	Yellow
Plasma Viscosity Level	Yellow	Yellow	Yellow	Yellow	Yellow	Yellow
Hemorheology	Yellow	Yellow	Yellow	Yellow	Red	Red
Decreased Vascular Flow	Red	Red	Red	Red	Red	Red
Hyperviscosity	Yellow	Yellow	Yellow	Yellow	Red	Green
Fibrinolysis	Green	Yellow	Yellow	Green	Yellow	Green
Thrombosis	Green	Yellow	Yellow	Green	Yellow	Green
Platelet Adhesiveness	Yellow	Yellow	Yellow	Yellow	Red	Red
Effects, blood coagulation	Yellow	Yellow	Yellow	Yellow	Red	Red
Vasodilatation	Yellow	Yellow	Yellow	Yellow	Red	Red
Vasodilation	Green	Green	Green	Green	Yellow	Green
Vasospasm	Green	Green	Yellow	Green	Yellow	Red
Vasospasm Mechanisms	Yellow	Red	Red	Yellow	Red	Red
Vasomotion	Yellow	Yellow	Yellow	Yellow	Yellow	Green
Decreased Vascular Resistance	Green	Green	Yellow	Yellow	Yellow	Red
<b>Total Found:</b>	<b>94.44%</b>	<b>88.88%</b>	<b>83.33%</b>	<b>94.44%</b>	<b>66.66%</b>	<b>61.11%</b>

Our results showed an ‘inverse effect’ between the open versus closed mode of discovery. In closed discovery a more specific, but smaller in size corpus (Raynaud disease) produced better results than a more general and bigger in size corpus (PVD or Hematology). On the contrary, in open discovery, a more general and bigger corpus (Vascular disease or PVD) produced better results than a more specific and smaller in size corpus (Raynaud disease).

### Conclusion

A novel LBD system utilizing UMLS ontologies, Semantic Filters and Word Embeddings was introduced which lead with minimal manual curation to the re-discovery of the Raynaud disease – Fish oil connection. Our LBD system offers an attractive alternative to other current drug repositioning approaches in order to accelerate the lead gathering process of finding a reliable vaccine to COVID-19 pandemic.

## References

- [1] Reed, T.S. (2020). Use of Word Embeddings in a Literature-Based Discovery System. Master of Science by Research Thesis. School of Computer Science, University of Lincoln, UK
- [2] Swanson, D. R. (1986). Fish oil, Raynaud's syndrome, and undiscovered public knowledge. *Perspectives in biology and medicine* 31(4): 7-18
- [3] Mikolov, T., Sutskever, I., Chen, K., Corrado, G., Dean, J. (2013). Efficient estimation of words and phrases and their compositionality. In *Proceedings of the 26th International Conference on Neural Information Processing Systems (NIPS 2013)*, 3111-3111. Lake Tahoe, Nevada
- [4] Aronson, A. A., Lang, F. M. (2010). An overview of Metamap: historical perspective and recent advances. *J Am Med Inform Assoc* 17: 229-236.
- [5] Bodenreider, O. (2004). The Unified Medical Language System (UMLS): integrating biomedical terminology. *Nuclei Acids Res* 32: D267-D270.