

Exploring Query Expansion for Entity Searches in PubMed

Chung-Chi Huang

National Center for Biotechnology Information (NCBI),
National Library of Medicine,
National Institutes of Health (NIH)

chuang@frostburg.edu

Zhiyong Lu

Zhiyong.lu@nih.gov

Abstract

Identifying relevant studies from the entire scientific literature is an important task in biomedical research. Past efforts have incorporated semantically recognized biological entities and medical ontologies into biomedical literature search. However, semantic relations are largely overlooked by biomedical search engines. In this work, we aim to discover synonymous biomedical semantic relations between entities and explore their uses in query (semantics) understanding for improved retrieval performance. Specifically, we discover synonymous semantic relations from PubMed queries and apply them to query expansion and specification. In these two real-world scenarios, better PubMed retrieval effectiveness, in terms of recall and precision, can be achieved, demonstrating the utility of our proposed approach.

1 Introduction

PubMed is widely used by millions of users on a daily basis for seeking scholarly publications in biology and life sciences. Recent studies show that a significant portion of PubMed queries are entity specific (i.e. entity searches) (Neveol et al., 2011; Huang and Lu, 2016).

Domain-specific search engines, such as PubMed, typically handle queries with domain knowledge in mind. For example, PubMed incorporates Medical Subject Headings (MeSH) to retrieve documents associated with query's semantic meaning than just keyword matches as in biomedicine it is common for concepts to appear in different forms in user queries and scholarly publications (Lu et al. 2009). However, PubMed can still suffer from mismatches between document and query words when an information need

involves entity semantic relations (Baumgartner et al., 2007).

Consider the query *chlorthalidone vs hydrochlorothiazide* and *chlorthalidone versus hydrochlorothiazide*. Semantically similar as they are, PubMed returns twice more relevant documents for the latter, clearly overlooking the semantics of the general terms of *vs* and *versus* during its search. Unfortunately, such performance difference resulting from different query formulations can lead to different levels of user satisfaction and different user experience with PubMed.

In light of this, we propose a framework where we first understand user query's semantics by discovering synonymous patterns among user queries (e.g. patterns *CHEMICAL vs CHEMICAL* and *CHEMICAL versus CHEMICAL*) for entity relations of interest. We then apply these learned synonymous patterns in query expansion to improve retrieval effectiveness for entity searches in PubMed.

In this work, we mine synonymous patterns in user queries instead of scholarly publications because queries are generally short (Islamaj Dogan et al., 2009; Wilkinson et al., 1995) and tend to bond entities in proximity. Here we specifically target chemical-chemical and chemical-disease relations such as chemical-induced-disease relation (Wei et al., 2016). The proposed framework, however, is easily generalizable to understand other bio-entity relations such as protein-protein interaction (Phizicky and Fields, 1995).

Our work is unique in several aspects. First, PubMed queries are semantically analyzed through context patterns, and synonymous relations or synonymous context patterns are discovered automatically. Second, synonymous patterns are applied to expand entity searches at pattern level to improve recall of relevant documents. Third, synonymous patterns can also be applied

to searches with entities only, where we add additional constraints to improve precision. Overall evaluation is able to point key directions for future development and improvement of PubMed, and can also shed light on how to effectively search biomedical literature beyond PubMed.

2 Related Work

Query Expansion (QE) has been an area of active research in Information Retrieval (IR). QE techniques manage to alleviate vocabulary mismatch between query and document words by adding related words to the initial queries, with the goal of improving retrieval effectiveness. Below we discuss three types of QE techniques classified based on how they derive related words: ontology-oriented, query-independent data-driven, and query-dependent data-driven technique.

Ontology-oriented techniques leverage language properties (e.g. synonyms, hypernyms and etc.) in dictionaries (Liddy and Myaeng, 1993), thesauri, or lexical databases (Voorhees, 1994) to find QE. General-purpose lexical database e.g. WordNet (Fellbaum, 1998) or a domain-specific one e.g. MeSH (Nelson et al., 2001) may be used.

Query-independent data-driven QE methods identify queries' similar words by analyzing global-wide documents not specific to queries. Hence, they are also known as global corpus-specific QE methods (Carpineto and Romano, 2012). They learn word association by concept terms (Qiu and Frei, 1993), term clustering (Crouch and Yang, 1992), distributional similarity (Lin 1998; Turney 2001; Chen et al., 2006), semantic topics (Park and Pamamohanarao, 2007), to name a few.

Query-dependent data-driven techniques, on the other hand, analyze query-specific documents for QE. While relevance feedback uses relevant documents from the initial queries, pseudo-relevance feedback uses top-ranked documents without human intervention (Xu and Croft, 1996). Measures for finding related terms in initially returned documents include Rocchio's weighting (Rocchio, 1971), Chi-square (Doszkocs, 1978), and Kullback-Leibler distance (Carpineto et al., 2001). Recently, Cui et al. (2003) and Riezler et al. (2007) consider user-clicked documents relevant for QE.

In biomedicine, QE studies primarily focus on ontologies and pseudo-relevance feedback. For example, Jalali and Borujerdi (2008) and Lu et al. (2009) expand queries via MeSH ontology,

and Srinivasan (1996), Aronson (1996), and Zhu et al. (2006) expand queries via Unified Medical Language System (Lindberg et al., 1993). On the other hand, biomedical queries can be reformulated (Lu et al., 2009) or systematically expanded based on initially retrieved documents focusing on abbreviations (Bacchin and Melucci, 2005), the controlled vocabulary of MeSH (Thesprasith and Jaruskulchai, 2014), or open vocabulary (Rivas et al., 2014).

In contrast to previous work, we semantically analyze frequently-sought general patterns (or relations) in biomedical queries, discover pattern synonyms, and use these automatically-learned synonymous patterns to expand real-world entity searches in PubMed. Such general-phrase pattern-level semantics understanding, complementary to domain-specific MeSH, later proves useful in QE and beneficial to PubMed literature search in our case studies.

3 Entity Searches in PubMed

(a) PubMed titles for the search *midazolam sevoflurane*

1. Network Meta-Analysis on the Efficacy of Dexmedetomidine, **Midazolam**, Ketamine, Propofol, and Fentanyl for the Prevention of **Sevoflurane**-Related Emergence Agitation in Children.
2. Determination of optimum time for intravenous cannulation after induction with **sevoflurane** and nitrous oxide in children premedicated with **midazolam**

(b) PubMed titles for its semantics-constrained query *midazolam vs sevoflurane OR midazolam versus sevoflurane OR ...*

1. Long-term sedation in intensive care unit: a randomized **comparison** between inhaled **sevoflurane** and intravenous propofol or **midazolam**.
2. Complications of **sevoflurane**-fentanyl **versus midazolam**-fentanyl anesthesia in pediatric cleft lip and palate surgery: a randomized comparison study.

Table 1. An example of PubMed search results (sorted by relevance) without (a) and with (b) semantic expansion.

We focus on understanding users' information needs or search semantics when they submit entity searches to PubMed. We discover synonymous patterns or entity relations in user queries (Section 3.1) and exploit them in the following two use scenarios to improve PubMed retrieval effectiveness.

Scenario 1. Consider an entity pair search with explicit relation mention (e.g. comparison relation between two drugs as in *albuterol vs levalbuterol*). We expand the query with its synonymous counterparts belonging to the same pattern-level relation (e.g. adding *albuterol versus levalbuterol*, *comparison between albuterol and levalbuterol*, and etc.). With such query expansion, we expect to retrieve

semantic relation	pattern	BiR	semantic relation	pattern	BiR
drug comparison	#C versus #C	2.38	drug-induced-disease	#C induced #D	1.14
	#C vs #C	10.05		#D induced by #C	1.14
	comparison of #C and #C	1.91		#D associate with #C	969.6
	comparison #C and #C	1.91		#C side effect #D	303
	#C compare #C	135.65		#D caused by #C	21.07
	difference between #C and #C	144		#C exposure and #D	21.26
	comparison between #C and #C	1.91		#C cause #D	48
	#C compare to #C	135.65		#D risk factor #C	94.13
	#C compare with #C	135.65		#D #C adverse effect	484.8
difference #C and #C	144				
drug combination	#C and #C combination	1.35	drug-treats-disease	#D treatment #C	1.96
	combination of #C and #C	1.35		#D and #C therapy	2.41
	combine #C and #C	904.2		treatment of #D with #C	1.96
	#C in combination with #C	1.35		treatment of #D #C	1.96
	#C and #C combination therapy	6.14		#D treatment with #C	1.96
	#C combined with #C	4.93		#C treatment for #D	1.96
	add #C to #C	37.99		#C in the treatment of #D	1.96
	combination therapy with #C and #C	6.14		#C in #D treatment	1.96
	concomitant #C and #C	38.64		#D treated with #C	7.59
			#C therapy in #D	2.41	

Table 2. Retrieval benefit in recall (BiR) when using synonymous relational patterns for query expansion.

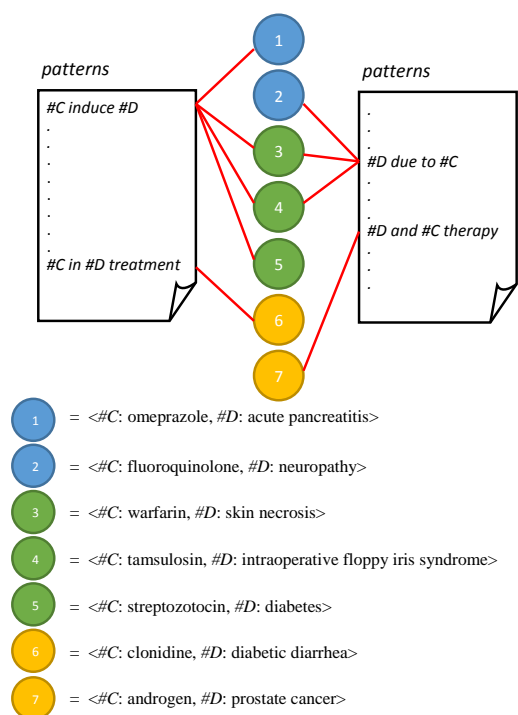


Figure 1. Patterns’ semantics similarity in terms of overlapping entities or LSA topics. While circles represent entities, the colors of the circles represent learned LSA topics.

from PubMed additional documents originally unreachable and expect to balance PubMed results across different query formulations with identical semantics meaning.

Scenario 2. Consider a pure entity pair search without any explicit mention of entity relation (e.g. *midazolam sevoflurane*). We constrain the query on its known search semantics learned based on

past PubMed searches (e.g. comparison relation between these two drugs). The newly constructed search (e.g. *midazolam vs sevoflurane OR midazolam versus sevoflurane OR ...* where OR combines PubMed results) is expected to direct PubMed towards documents users truly interested in but otherwise might be ranked low based on the original search. Take Table 1 for example. Top-ranked documents are more relevant with the new semantics-constrained query if users are to compare the two entities without explicitly mentioning so in the search query.

Note that in *Scenario 1* and *Scenario 2* we perform PubMed searches under relevance sorting (as opposed to the default chronological sorting) and we search PubMed and use matches in article titles as a proxy for human relevance evaluation (Kim et al., 2016). In other words, to ensure quick turnaround and large-scale evaluation, we assume those matching titles all satisfy users’ information needs (i.e. perfect precision) and thus no human relevance judgments is required.

3.1 Discovering Synonymous Patterns

We have previously developed an unsupervised approach for identifying synonymous patterns of entity relations in PubMed queries (Huang and Lu, 2016). Due to space limitation, we only briefly outline major steps below. We refer interested readers to (Huang and Lu, 2016) for details.

First, a six-month worth of PubMed queries (35M queries) are stemmed and tagged using entity recognition tools (Wei et al., 2015;

Leaman et al., 2013; Leaman et al., 2015) for genes/proteins, diseases, and chemicals/drugs.

Next, we formulate queries to context patterns and focus on specifically discovering synonymous patterns for chemical-chemical (*CC*) and chemical-disease (*CD*) relations. For instance, the query *skin necrosis associate with warfarin* is formulated into *#D associate with #C* where *#C* and *#D* stands for chemical and disease entity respectively.

Inspired by distributional similarity (Lin 1998), we then exploit these patterns' participating entity pairs to understand their semantics. In such a way, synonymous patterns can be found in an unsupervised fashion in contrast to seeds-required pattern recognition work (e.g. Xu and Wang, 2014). Take Figure 1 for example. Our framework will consider the pattern *#C induce #D* semantically closer to *#D due to #C* than to *#C in #D treatment since #C induce #D* and *#D due to #C* share more participating entities in user queries: 2 overlapping entities out of 7 entities vs 0 out of 7.

To avoid data sparseness issue on (distributional similarity in) entity mention, we further leverage latent semantic analysis, LSA, (Rehurek and Sojka, 2010) to find entities' LSA topics which in turn reduces the space of semantics analysis from the dimension of entity pairs to a much smaller dimension of LSA topics. The benefit of using LSA topics is clear: after LSA transformation, *#C induce #D* in Figure 1, where circle's colors depict LSA topics, shows stronger semantics connection with *#D due to #C* than previously without LSA: 2 overlapping LSA topics out of 3 topics.

Our LSA-based approach is able to achieve satisfying performance in finding semantically similar patterns across entity relations of interest, such as drug-induced-disease relation, drug-drug interaction, to name a few. We refer interested readers to (Huang and Lu, 2016) for detailed evaluation results.

3.2 Expanding Entity-Relation Searches

Once our method identifies candidates of pattern synonyms, we collect the set of true synonymous patterns and apply them to semantic query expansion as below.

We first order a semantic relation's synonymous patterns according to their frequencies in PubMed queries, which represent user preferences or user intuitions (in searching the target bio-relation between two entities). See patterns in descending order of frequency in the second

and fifth column of Table 2. For example, PubMed users prefer using *#C versus #C* to *#C vs #C* or *comparison of #C and #C* in comparing two drugs. Currently, four common entity relations between drugs and between drugs and diseases are of our particular interest: drug comparison, drug combination, drug-induced-disease and drug-treats-disease.

Second, for each relation, we assemble its 500 most searched entity pairs from our search logs. For example, *<albuterol, levalbuterol>* is a popular chemical pair for the drug comparison relation.

For each entity pair (e.g. *<albuterol, levalbuterol>*) of a semantic relation, we then submit a query with the pair using one of the relational patterns (e.g. *albuterol vs levalbuterol*) and compare the search result with that of semantically expanded query that leverages all synonymous patterns (e.g. *albuterol versus levalbuterol OR albuterol vs levalbuterol OR ...* Syntax OR combines PubMed retrieval results). Recall that the searches are limited to PubMed titles. Finally, we compute the ratio of the number of total search results via all patterns of the semantic relation over that of each individual pattern, averaged over 500 entity pairs. Such difference in recall is referred to as benefit in recall, **BiR**.

As Table 2 shows, a **BiR** score above 1 means expanding queries using collective synonymous patterns of the same semantics improves PubMed recall or helps PubMed retrieve more relevant documents. Take the drug comparison relation for example. Regardless of the chemical pair of interest, expanded queries can always retrieve more relevant documents than using the individual pattern of *#C versus #C* (more than twice as many on average: 2.38). In some cases of Table 2, the improvement in recall is substantial (e.g. 135.65 associated with *#C compare #C*, 904.2 associated with *combine #C and #C*, and so on).

The benefit of using our synonymous patterns for query expansion in current PubMed settings can be observed across various types of *CC* or *CD* entity-relation searches, searches with explicit relation mention. And interestingly, the most frequently used patterns by users (or the most intuitive/straightforward search patterns from users' points of view) may not always be the best choice at default: among the drug comparison patterns, *comparison of #C and #C* is more effective than the most popular *#C versus #C* in retrieving relevant documents. A semantic framework like ours can balance PubMed retrieval results across different entity-relation expressions in searches with similar meanings.

3.3 Expanding Pure Entity Pair Searches

Among PubMed searches, pure entity pair searches or searches containing only two bio-entities without any explicit relation mentions (e.g. *midazolam sevoflurane*), account for approximately half of the searches involving dual bio-entities. As a result, we investigate in this subsection how we can improve PubMed user experience by expanding these queries, with the help of our synonymous patterns and past user searches. The process is detailed below.

First, we identify pure entity pair searches only sought by PubMed users in *a specific* relation/context, based on which we expand the searches and impose semantic search constraints. Take the pure entity pair search *midazolam sevoflurane* for instance. Since it had only been searched with drug comparison relation by PubMed users, we later explicitly constrain that search query in the context of drug comparison relation. This step infers the implicit relation between the entity pair from the wisdom of the crowd (i.e. past search logs). Our hypothesis is that such implicit relation, if explicitly added to the search, may improve retrieval results and in turn user experience.

In the current experiment, a total of 1,600 unique pure entity-pair queries are collected with *CC* relation constraints (i.e. drug comparison, drug combination, and drug interaction) and *CD* relation constraints (i.e. drug-treats-disease, drug-induced-disease, supplement-for-disease, drug-resistance-in-disease).

Similar to the settings in Section 3.2, we submit to PubMed (a) original queries, i.e. pure entity pairs and (b) expanded queries with explicit relation constraints learnt from past user queries. For example, original search *midazolam sevoflurane* and its semantics-constrained counterpart *midazolam versus sevoflurane OR midazolam vs sevoflurane OR ...* (expanded using our synonymous patterns of the drug comparison relation, in which *midazolam sevoflurane* had only been sought) will be submitted to PubMed.

Finally, based on the search results from (a) and (b), we compute the retrieval effectiveness of regular PubMed by using (b)'s results as the ground truth. In other words, we assume the expanded queries truly represent users' search intention and their search results truly satisfy users' information needs. Retrieval performance is measured by standard information retrieval (IR) measures: precision (P), mean reciprocal rank

(MRR) and nDCG (Jarvelin and Kekalainen, 2002) at rank 20.

As we can see in Table 3, the difference between current performance scores in MRR or nDCG and perfect scores (i.e. perfect MRR or nDCG equals 1) suggests genuinely there is room for performance increase in retrieval for such searches, i.e. pure entity pair searches, in current PubMed settings. While pure *CD* searches yield better results than pure *CC* searches, potential gain in performance is still substantial for *CD* queries, which can be achieved by simply adding semantics constraints and expanding queries accordingly. In some cases (e.g. pure entity pair searches with implicit drug interaction relation), semantics constraints almost warrant a more satisfying search performance.

entity pair type	implicit relation	IR measures @ 20	results
CC	drug comparison	P	0.25
		MRR	0.43
		nDCG	0.57
	drug combination	P	0.29
		MRR	0.47
		nDCG	0.61
drug interaction	P	0.13	
	MRR	0.32	
	nDCG	0.43	
CD	drug-treats-disease	P	0.34
		MRR	0.58
		nDCG	0.66
	drug-induced-disease	P	0.36
		MRR	0.63
		nDCG	0.70
	supplement-for-disease	P	0.23
		MRR	0.47
drug-resistance-in-disease	nDCG	0.56	
	P	0.21	
	MRR	0.43	
	nDCG	0.55	

Table 3. Results on pure *CC* and *CD* queries with implicit relations.

4 Summary

We have applied query semantics understanding to PubMed literature search. The proposed framework involves discovering synonymous relational patterns in queries and, based on those, expanding PubMed user queries, specifically entity search queries. Preliminary evaluation shows such semantic query expansion helps to improve PubMed retrieval effectiveness. And better PubMed performance implies better user experience and less curation effort (Lu and Hirschman, 2012). Incorporating such general-phrase semantics framework, complementary to domain-specific MeSH, into PubMed serving millions of users is warranted.

5 Acknowledgements

This work was supported by the Intramural Research Program of the National Library of Medicine, National Institutes of Health. The authors would like to thank anonymous reviewers for their suggestions and comments.

Reference

- Aronson, A.R. 1996. The effect of textual variation on concept based information retrieval. *Proc AMIA Annu Fall Symp.*
- Aronson, A.R. and T.C. Rindfleisch. 1997. Query expansion using the UMLS Metathesaurus. *Proc AMIA Annu Fall Symp.*
- Bacchin, M. and M. Melucci. 2005. Symbol-based query expansion experiments at TREC 2005 Genomics Track. In *Proceedings of Text REtrieval Conference.*
- Baumgartner, W, Z. Lu, H. Johnson et al. 2007. An integrated approach to concept recognition in biomedical text. In *Proceedings of BioCreative Challenge Evaluation Workshop.*
- Carpineto, C., R. De Mori, G. Romano et al. 2001. An information-theoretic approach to automatic query expansion. *ACM Transactions on Information Systems.*
- Carpineto, C. and G. Romano. 2012. A survey of automatic query expansion in information retrieval. *ACM Computing Surveys.*
- Chen, H., M. Lin, and Y. Wei. 2006. Novel association measures using web search with double checking. In *Proceedings of ACL*, p. 1009-1016.
- Crouch, C.J. and B. Yang. 1992. Experiments in automatic statistical thesaurus construction. In *Proceedings of ACM SIGIR.*
- Cui, H., J.R. Wen, J.Y. Nie et al. 2003. Query expansion by mining user logs. *IEEE Transactions on Knowledge and Data Engineering.*
- Deerwester, S., S.T. Dumais, G.W. Furnas et al. 1990. Indexing by latent semantic analysis. *Journal of the Association for Information Science.*
- Diaz-Galiano, M.C., M.T. Martin-Valdivia, and L.A. Urena-Lopez. 2009. Query expansion with a medical ontology to improve a multimodal information retrieval system. *Comput Biol Med.*
- Dramé, K., F. Mougin, and G. Diallo. 2014. Query expansion using external resources for improving information retrieval in the biomedical domain. In *Proceedings of ShARe/CLEF eHealth Evaluation Lab.*
- Doszkocs, T.E. 1978. AID, an associative interactive dictionary for online searching. *Online Review.*
- Fellbaum, C. 1998. WordNet: an electronic lexical database.
- Gauch, S., J. Wang, and S.M. Rachakonda. 1999. A corpus analysis approach for automatic query expansion and its extension to multiple Databases. *ACM Transactions on Information Systems.*
- Gonzalo, J., F. Verdejo, I. Chugur et al. 1998. Indexing with WordNet synsets can improve text retrieval. In *Proceedings of ACL Workshop.*
- Huang, C.C. and Z. Lu. 2016. Discovering biomedical semantic relations in PubMed queries for information retrieval and database curation. *Database.*
- Islamaj Dogan, R., G.C. Murray, A. Neveol et al. 2009. Understanding PubMed user search behavior through log analysis. *Database.*
- Jalali, V. and M.R.M. Borujerdi. 2008. The effect of using domain specific ontologies in query expansion in medical field. In *Proceedings of IEEE International Conference on Innovations in Information Technology.*
- Jarvelin, K. and J. Kekalainen. 2002. Cumulated gain-based evaluation of IR technologies. *ACM Transactions on Information Systems.*
- Kim, S., W.J. Wilbur, Z. Lu. 2016. Bridging the gap: a semantic similarity measure between queries and documents. arXiv:1608.01972.
- Lavrenko, V. and W.B. Croft. 2001 Relevance based language models. In *Proceedings of ACM SIGIR.*
- Leaman, R., R. Islamaj Dogan, and Z. Lu. 2013. DNorm: disease name normalization with pairwise learning to rank. *Bioinformatics.*
- Leaman, R., C.H. Wei, and Z. Lu. 2015. tmChem: a high performance approach for chemical named entity recognition and normalization. *J Cheminform.*
- Liddy, E.D. and S.H. Myaeng. 1993. DR-LINK's linguistic-conceptual approach to document detection. In *Proceedings of Text REtrieval Conference.*
- Lin, D. 1998. Automatic retrieval and clustering of similar words. In *Proceedings of ACL*, p. 768-774.
- Lindberg, D.A., B.L. Humphreys, and A.T. McCray. 1993. The Unified Medical Language System. *Methods Inf Med.*
- Lu, Z. and L. Hirschman. 2012. Biocuration workflows and text mining: overview of the BioCreative 2012 Workshop Track II. *Database.*
- Lu, Z., W. Kim, and W.J. Wilbur. 2009. Evaluation of Query Expansion Using MeSH in PubMed. *Inf Retr.*

- Lu, Z., W.J. Wilbur, J.R. McEntyre et al. 2009. Finding query suggestions for PubMed. In *AMIA Annu Symp Proc*.
- Nelson, S.J., W.D. Johnston, and B.L. Humphreys. 2001. Relationships in medical subject headings (MeSH).
- Neveol, A., R. Islamaj Dogan, and Z. Lu. 2011. Semi-automatic semantic annotation of PubMed queries: a study on quality, efficiency, satisfaction. *J Biomed Inform*.
- Park, L.A.F. and K. Ramamohanarao. 2007. Query expansion using a collection dependent probabilistic latent semantic thesaurus. In *Proceedings of PAKDD*.
- Phizicky, E.M. and S. Fields. 1995. Protein-protein interactions: methods for detection and analysis. *Microbiol Rev*.
- Qiu, Y. and H.P. Frei. 1993. Concept based query expansion. In *Proceedings of ACM SIGIR*.
- Rehurek, R. and P. Sojka. 2010. Software framework for topic modelling with large corpora. In *Proceedings of LREC Workshop*.
- Riezler, S., E. Vasserman, I. Tsochantaridis et al. 2007. Statistical machine translation for query expansion in answer retrieval. In *Proceedings of ACL*.
- Rocchio, J.J. 1971. Relevance feedback in information retrieval.
- Srinivasan, P. 1996. Query expansion and MEDLINE. *Information Processing and Management*.
- Thesprasith, O. and C. Jaruskulchai. 2014. Query expansion using medical subject headings terms in the biomedical documents. *Intelligent Information and Database Systems*.
- Tsuruoka, Y. and J. Tsujii. 2005. Bidirectional inference with the easiest-first strategy for tagging sequence data. In *Proceedings of EMNLP*, p. 467-474.
- Turney, P.D. 2001. Mining the Web for synonyms: PMI-IR versus LSA on TOEFL. In *Proceedings of EMCL*, p. 491-502.
- Voorhees, E.M. 1994. Query expansion using lexical-semantic relations. In *Proceedings of ACM SIGIR*.
- Wei, C.H., H.Y. Kao, and Z. Lu. 2015. GNormPlus: an integrative approach for tagging genes, gene families, and protein domains. *Biomed Res Int*.
- Wei, C.H., Y. Peng, R. Leaman et al. 2016. Assessing the state of the art in biomedical relation extraction: overview of the BioCreative V chemical-disease relation (CDR) task. *Database*.
- Wilkinson, R., J. Zobel, and R. Sacks-Davis. 1995. Similarity measures for short queries. In *Proceedings of Text REtrieval Conference*.
- Xu, R. and Q. Wang. 2014. Automatic construction of a large-scale and accurate drug-side-effect association knowledge base from biomedical literature. *J Biomed Inform*.
- Xu, X. and X. Hu. 2010. Cluster-based query expansion using language modeling in the biomedical domain. In *Proceedings of IEEE International Conference on Bioinformatics and Biomedicine Workshops*.
- Zhai, C. and J. Lafferty. 2001. Model-based feedback in the language modeling approach to information retrieval. In *Proceedings of CIKM*.
- Zhu, W., X. Xu, X. Hu et al. 2006. Using UMLS-based re-weighting terms as a query expansion strategy. In *Proceedings of IEEE International Conference on Granular Computing*.