# Application of Semantic Knowledge Representation and Natural Language Processing to Identify Pharmacologic Mechanisms

Sanya Bathla Taneja

*University of Pittsburgh, Pittsburgh, PA, USA*

### Abstract

Natural product-drug interactions (NPDIs) occur due to co-consumption of drugs and natural products leading to therapeutic failure and adverse events. Understanding the pharmacologic mechanisms of interaction is key to prevent adverse effects and improve drug safety. Major challenges in identification of NPDI mechanisms include variability in natural product composition and constituents, limited known pharmacokinetic information about constituents, and unavailability of gold standard datasets for NPDI mechanisms. I hypothesize that a large-scale, heterogeneous, biomedical knowledge graph (KG) combining biomedical ontologies, drug databases, and domain-specific scientific literature will represent relationships of natural products with other biomedical entities and will generate biologically plausible mechanisms for scientific research. In this work, I will construct a natural products-relevant KG and apply discovery methods, including discovery patterns, graph algorithms, and translational embedding methods to generate mechanistic hypotheses for 30 selected natural products. Mechanism generation in the KG will be guided by pharmacovigilance signals from spontaneous reporting systems. The evaluation will focus on (a) prediction of NPDIs and mechanisms using a reference dataset and (b) a user study to evaluate the quality of evidence for NPDIs in the KG to identify gaps for further research.

### Keywords

Knowledge graph, Biomedical ontology, Literature-based discovery, Natural products [1]

## 1. Introduction

The World Health Organization (WHO) estimates that up to four billion people use medicinal plants as healthcare and that the concomitant use of complementary health approaches and pharmaceutical drugs is widespread across the world[1]. Approximately 50% of adults in midlife have reported co-consumption of natural products and drugs, with the prevalence being even higher in older adults (up to 88%) in the United States (US) [2,3]. Such concomitant use of natural products and drugs can result in natural product-drug interactions (NPDIs) leading to therapeutic failure or adverse events[4,5]. Over the years, numerous studies have explored computational methods for mechanism discovery of pharmaceutical drugs including drug-target predictions, drug repurposing, and drug-drug interactions, and a large number of these studies have used knowledge graphs for prediction[6]. However, computational research on botanical and other natural products used for complementary health is not as widespread. Understanding the biochemical mechanisms underlying clinically significant NPDIs can help prevent or minimize adverse drug reactions (ADRs) resulting from the NPDIs[4].

Major challenges in the discovery of NPDI mechanisms include the variability in natural product compositions, challenges in identification of causative constituents, and limited pharmacokinetic information about the known constituents[1,4]. Pharmacokinetic NPDIs occur if a natural product extract (e.g., some quantity of green tea) phytoconstituent (e.g., catechin) inhibits or induces the function of a drug metabolizing enzyme or transporter, which may or may not have unforeseen negative consequences. Systematic literature reviews are used by

researchers to understand the research gaps, select natural products for further investigation, and design studies. The mechanistic hypotheses suggested in the literature inform the design of future experiments. Evaluating each natural product-drug pair on the market for a potential NPDI is thus time-consuming and expensive. Although progress has been made recently to overcome challenges[5,7], the increasing sales of natural products in the market, changing regulatory landscape, and growing safety concerns over NPDIs call for novel methods to help scientists make accurate and timely NPDI predictions. A biomedical knowledge graph (KG) combines expert-derived information sources into a graph where the nodes represent biomedical entities and edges represent relationships between the entities[6]. When integrated with domain-specific scientific literature through semantic relation extraction and named entity recognition methods, a KG is a powerful tool that can be used for cost-effective prediction and mechanism identification for NPDIs to guide researchers to identify research gaps and prioritize new experiments.

In this research, I propose to apply knowledge representation and natural language processing methods to construct a natural products-relevant KG that combines existing biomedical knowledge through ontologies with literature-based discovery. I will construct a semantically integrated KG that combines biomedical ontologies with full texts of domain-specific scientific literature. Plausible mechanistic hypotheses for potential NPDIs and associated ADRs will be generated using computational discovery methods such as discovery patterns and presented to researchers. The generation of biologically plausible mechanisms will be guided by pharmacovigilance signals from natural product spontaneous reports in the US Food and Drug Administration Adverse Event Reporting System (FAERS) and published case reports related to NPDIs. As there does not exist a gold standard dataset for NPDI mechanisms, I will also create a reference dataset for selected natural products to evaluate the KG. The quality of evidence available in the KG for NPDIs will then be evaluated with a pilot user study that presents NPDI mechanisms to researchers.

## 2. Related Work

Existing computational methods that predict NPDIs have focused on classification of NPDIs from literature or existing databases[8,9], designing literature retrieval systems[10], and relation extraction for herb-drug interactions[11]. Classification of NPDIs has been done using scientific abstracts, existing NPDI databases, and transfer learning approaches. The major challenges in computational discovery of NPDIs involve a lack of gold standard data on NPDIs and difficulty in obtaining representations of natural products and their constituents[8]. To overcome these challenges, studies have extracted knowledge from scientific abstracts from PubMed[8,12,13], used existing databases for reference data[8,14], and trained models using drug-drug interaction data[9]. Advanced methods such as graph representation learning have also been applied to classify food-drug interactions, although the results on external datasets have proved modest[8]. Follow-up evaluations of these methods are also lacking, and besides a study by Schutte et. al.[12], none of the studies have tried to elucidate the mechanisms underlying the interactions. Existing dietary supplement information retrieval tools targeted at consumers have successfully created graph-based visualizations with interactive features to provide information regarding dietary supplement uses, interactions, and ingredients[15,16].

As existing research in computational discovery of NPDIs using artificial intelligence methods has focused broadly on classification or retrieval of the interactions only, there exists a gap in research methods that can generate explainable mechanistic hypotheses for potential NPDIs and associated ADRs for scientists. Discovery patterns are interpretable sequences of nodes and relations in KG and have been successfully applied in literature-derived KGs to identify mechanistic information[17]. Recent graph representation learning methods that generate embeddings from the KG have also shown promise in discovering new edges in KGs and identifying mechanisms based on the similarity of nodes and edges in the embedding space[18]. Understanding the mechanistic explanations and available evidence for the mechanisms is particularly crucial for NPDI researchers to design and prioritize new studies. Using a

comprehensive representation of natural products-relevant knowledge with other biomedical entities through the incorporation of data from various sources and leveraging discovery patterns and embedding methods within the KG, we can identify the underlying mechanisms that for NPDIs as well as present them to researchers with corresponding metadata and supporting evidence.

To this end, this study will produce the first ontology-grounded KG focused on NPDIs that integrates heterogeneous data sources, including biomedical ontologies, open databases, and full texts of domain-specific scientific literature. The integration of sources ensures that the generated mechanisms are grounded in published results and existing biological knowledge. This study will be the first to integrate full texts of NPDI literature with the ontology-grounded KG using two high performance relation extraction systems. Unlike most literature-derived KGs that use scientific abstracts to extract information, using the full texts of literature to create the literature graphs derives mechanistic information that is not always present in the abstracts. This is also the first study to use pharmacovigilance signals from spontaneous reporting systems to guide the mechanism discovery for NPDIs and thus is also able to focus on the outcomes of the potential NPDIs. Finally, the evaluation strategies will be used to assess the plausibility of the mechanisms and reliability of results. Overall, the research presents a significant step forward in computational discovery of NPDIs which have the potential to improve drug safety and clinical decision making.

## 3. Research Hypotheses

**Hypothesis 1**: The integration of a large-scale, ontology-grounded KG with domain-specific scientific literature will provide an interconnected representation of natural products with other biomedical entities for recapturing knowledge about the natural products.
**Strategy:** To test this hypothesis, I will first construct a KG for natural products combining biomedical ontologies, databases, and domain-specific literature for 30 selected natural products. The selected natural products will be a mix of well- and less-known products. A natural language processing pipeline with semantic relation extraction and named entity recognition will be used to construct literature-based graphs for the natural products for integration in the KG. Metadata from all data sources and supporting data from the literature will be included in the KG to support the mechanism discovery and provide evidence for each link in the KG. To evaluate the KG, I will use discovery patterns and shortest path searches to recapture pharmacokinetic knowledge in the KG for two model natural products, green tea and kratom and their interacting enzymes and transporters and compare the information with human-curated data from the Center of Excellence for Natural Product Drug Interaction Research (NaPDI Center) database[5].

**Hypothesis 2**: Combining discovery patterns, graph algorithms, and embeddings in the KG will help to generate biologically plausible mechanistic hypotheses for potential NPDIs and pharmacovigilance signals of natural product-ADRs. The proposed KG will provide improved support to researchers in identifying research gaps for the natural product-related interactions when compared to existing approaches.
**Strategy:** I will create a reference dataset of known NPDIs and associated ADRs for the 30 natural products from existing resources, including the NaPDI database, Stockley's Herbal Medicines Interactions[19], and the Food Interactions with Drugs Evidence Ontology[20]. Then, I will apply discovery methods in the KG to predict NPDIs and associated ADRs, provide plausible mechanistic hypotheses, and compare the results with the reference dataset. The discovery methods will include discovery patterns, shortest path searches, and translational graph embedding methods to (a) predict the interaction between and (b) generate biologically plausible mechanistic hypotheses for an input natural product-drug or natural product-ADR pair. Mechanisms from the embeddings will be generated based on the cosine similarity between node vectors in the KG between the input node pairs. The focus will be on identifying mechanistic explanations for potential NPDIs and associated ADRs with reported signals in the FAERS

database. For evaluation, I will first calculate the accuracy, precision, and recall of predictions from the KG when compared to the reference dataset for NPDIs with known mechanisms. Then, I will create a user interface that displays the mechanisms to researchers with supporting data and evaluate the quality of evidence in the KG and usefulness in identifying research gaps for the selected natural products through a within-subjects user study. Further details of evaluation are in Section 5.
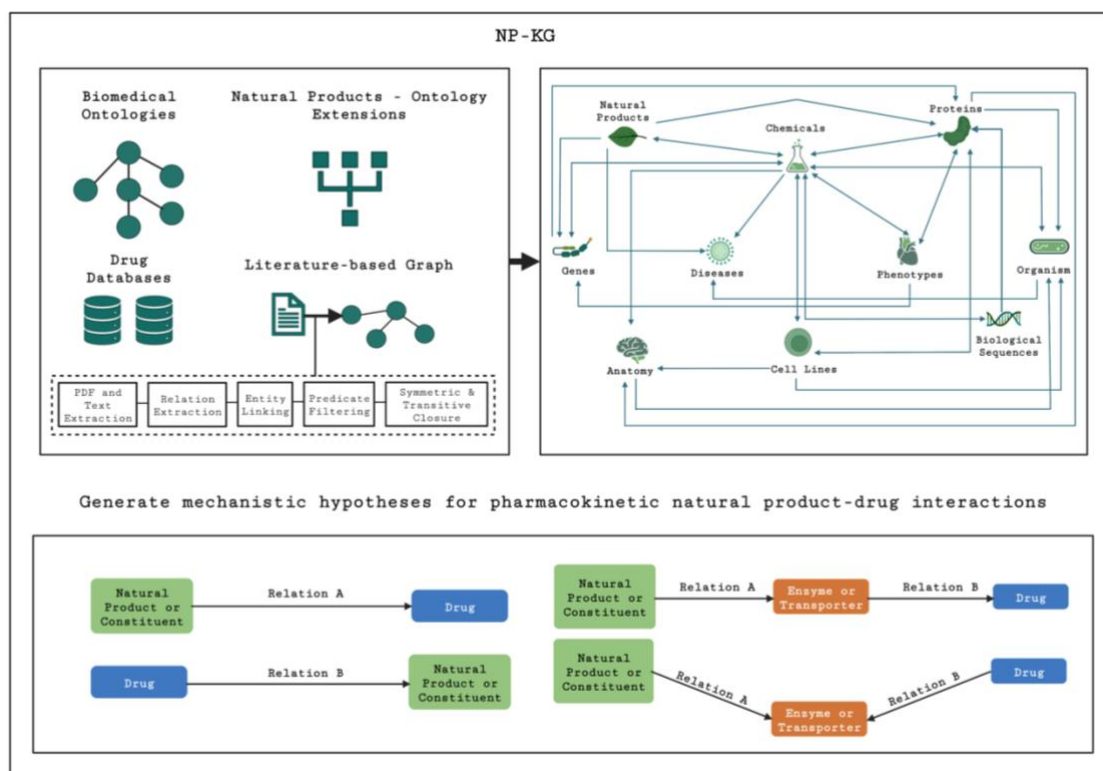
## 4. Preliminary Results



**Figure 1**: Overview of preliminary work to construct the KG and apply discovery patterns for NPDIs with interacting enzymes and transporters.

In preliminary work, I developed a heterogeneous KG for 30 natural products that combined an ontology-grounded KG with a literature-based KG (Figure 1). The natural products were selected based on published case reports analysis of NPDIs and pharmacovigilance signals from natural products adverse event reports from the FAERS database. The ontology-grounded KG was constructed using the PheKnowLator workflow that semantically integrates 13 Open Biological and Biomedical Ontologies (OBO) Foundry ontologies and linked data sources for biomedical entities such as chemicals, proteins, diseases, genes, phenotypes, and more[21,22]. Natural products were included as ontology extensions in the ChEBI lite ontology and integrated in the KG[23].

The literature-based KG was constructed from relation extraction of full texts of scientific publications for the natural products after applying search strategies in PubMed. Predications were extracted from two relation extraction systems, SemRep[24] and the Reading and Assembling Contextual and Holistic Mechanisms from Text (REACH) biological reader with the Integrated Network and Dynamic Reasoning Assembler (INDRA) framework[25,26]. The scope of the literature for the literature-based KG included all PubMed-indexed articles related to natural products (including keywords of scientific names, synonyms, and their constituents) and pharmacokinetic interactions. The ontology-grounded and literature-based KGs were semantically integrated after linking all subjects, predicates, and objects from the literature-

based KG to OBO concepts using both manual and automated entity linking methods. The combined KG is termed NP-KG[27].

The combined NP-KG contained 1,090,172 nodes and 7,920,893 edges. It is publicly available in both serialized and gpickle formats[28]. The ontology-grounded KG contained 1,089,613 nodes and 7,836,662 edges. The literature-based graph constructed from the combined and deduplicated predications of natural products contained 8,782 nodes and 84,569 edges. The literature-based graph added 559 unique nodes and 84,231 unique edges from 3,508 full texts processed by SemRep and 4,318 full texts processed by REACH. NP-KG contains all relevant metadata from databases and publications, including year, source, source sentence, study type, source section of publication, source sentence, and reference as edge attributes in the KG. Semantic representations were created for 30 natural products and 571 unique constituents. Out of the 571 unique constituents, 153 (26.8%) did not already exist in ChEBI ontology and were added as new classes. After integrating the natural products and constituents, 255 classes and 3695 axioms were added to ChEBI Lite ontology, bringing the total to 182,629 classes and 1,398,337 logical axioms.

The evaluation strategy included knowledge recapturing through shortest path searches and application of discovery patterns for two model natural products, green tea and kratom to find interacting enzymes, transporters, and NPDI mechanisms in the KG when compared to human curated data from the NaPDI Center database[5]. The evaluation aimed to recapture known information about interacting enzymes and transporters for green tea- and kratom-related pharmacokinetic NPDIs in NP-KG and establish congruence or contradiction when compared to ground truth information. Table 1 summarizes the results of direct edges and shortest path searches for congruent and contradictory information in the KG. For the green tea-related nodes, I performed 59 searches for direct edges or shortest paths in NP-KG involving 19 enzymes and 8 transporters (39.98% congruent, 15.25% contradictory, 3.39% both). For the kratom-related nodes, I performed 14 searches for direct edges or shortest paths involving 10 enzymes and 1 transporter (50% congruent, 21.43% contradictory, 7.14% both). Results with both congruent and contradictory edges between the nodes were manually reviewed to verify congruence and/or contradiction and for error analysis. Further, discovery patterns shown in Figure 1 were applied for five natural product-drug pairs, including green tea-nadolol, green tea-raloxifene, kratom-midazolam, kratom-quetiapine, and kratom-venlafaxine, with known interactions to find hypotheses for potential pharmacokinetic NPDIs. The preliminary results showed that the KG can capture information about the interacting enzymes and transporters and generate mechanistic hypotheses for the interacting natural product-drug pairs. The KG further successfully identified interacting enzymes and transporters for the natural product-drug pairs as shown in Figure 2.

**Table 1**
**Summary of congruences and contradictions for direct edges and shortest paths in the KG compared to ground truth information for green tea and kratom.**

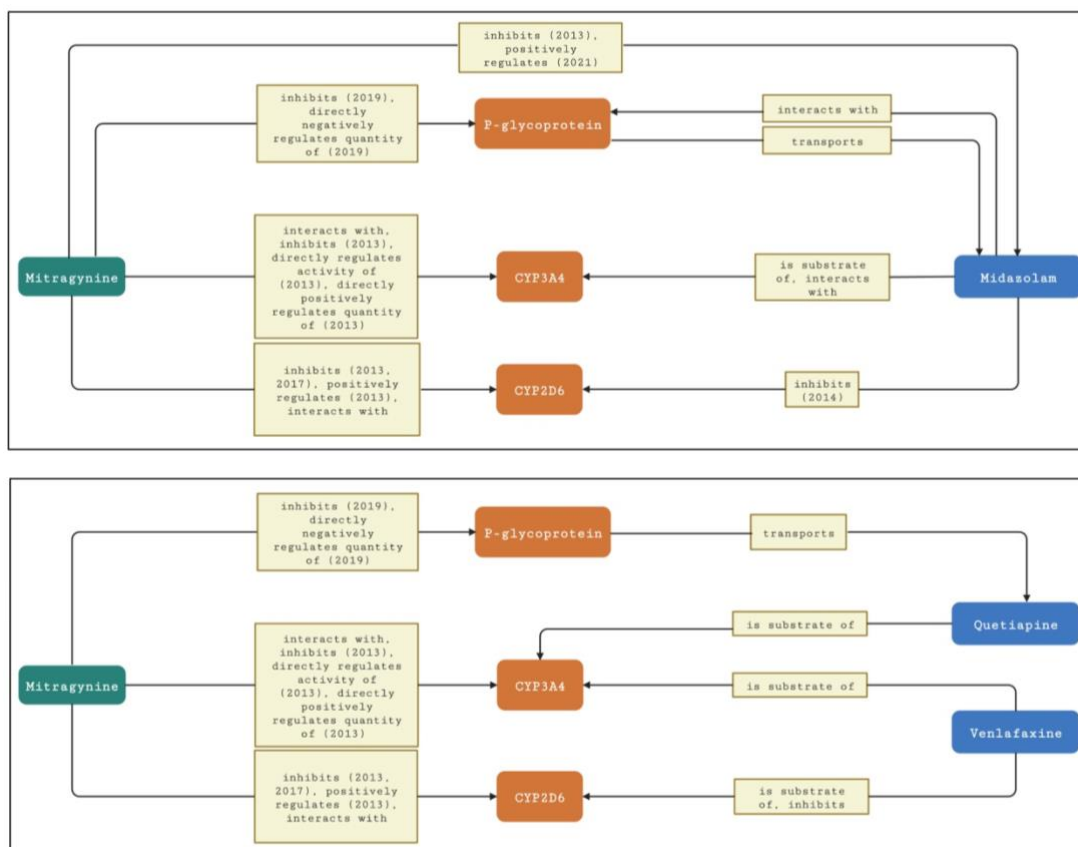| Head 1 | Green Tea (%) | Kratom (%) |
|---|---|---|
| Congruence | 23 (38.98) | 7 (50.0) |
| Contradiction | 9 (15.25) | 3 (21.43) |
| Edges/paths exist but no congruence or contradiction | 25 (42.37) | 3 (21.43) |
| Both congruence and contradiction | 2 (3.39) | 1 (7.14) |
| Total searches | 59 | 14 |

**Figure 2**: Discovery pattern results for kratom-midazolam, kratom-quetiapine, and kratom-venlafaxine with interacting enzymes (Cytochrome P450 (CYP) 3A4 and 2D6) and a transporter (P-glycoprotein). Rounded rectangles represent nodes and rectangles represent edges in the KG. If an edge is derived from the literature, the year of publication is noted with the edge label.

## 5. Evaluation

The preliminary work evaluated the potential of the KG to recapture known information about the interacting enzymes and transporters for two model natural products, green tea and kratom when compared to human curated data. In future work, I will scale the methods to evaluate the performance of the discovery methods, including discovery patterns, graph algorithms, and translational embedding methods for known interactions in a reference dataset constructed from NPDI mechanistic data. Performance metrics, including accuracy, precision, and recall will be calculated and the plausibility of the mechanisms will be evaluated through a review by pharmacists. The metrics will be calculated for two versions of the KG, a time-sliced version with data from 2021 and prior, and a version with data from 2023 to evaluate the ability of the KG to discover new knowledge.

Then, I will create a prototype tool that presents the mechanisms to researchers along with metadata and supporting data (including publication details such as year and study type, measurement information, source of data) for each link in the mechanism for potential NPDIs reported in the FAERS database. The quality of evidence available in the KG will be evaluated through a within-subjects user study that compares the use of existing methods (literature review) and KG-generated mechanisms for identifying research gaps for NPDIs, with the hypothesis that the proposed KG and discovery methods will better support researchers in identifying research gaps for NPDIs and lead to quicker resolution of questions when compared to existing approaches. This will be tested based on a questionnaire designed by NPDI experts with NPDI-related questions and semi-structured interviews with the participants of the user study.

# 6. Conclusion

The proposed work presents methods for construction of a large-scale biomedical KG combining biomedical ontologies, drug databases and full texts of domain-specific scientific literature to generate mechanistic hypotheses for NPDIs. Preliminary work has shown the potential of the KG to capture known mechanistic information about two model natural products and interacting enzymes and transporters. More advanced discovery methods, including translational embedding methods can now be used to predict NPDIs and generate mechanistic hypotheses from the KG using maximization of cosine similarity and path degree product of the embedding vectors. Then a combination of the discovery methods can be applied to the KG to produce plausible mechanisms. The next steps will be to apply advanced discovery methods in the KG for a wider set of natural products and associated ADRs and evaluate the plausibility of the mechanisms and usefulness for researchers.

## Acknowledgements

## References

[1] Organization WH, others. Key technical issues of herbal medicines with reference to interaction with other medicines 2021.

[2] Agbabiaka TB, Wider B, Watson LK, Goodman C. Concurrent Use of Prescription Drugs and Herbal Medicinal Products in Older Adults: A Systematic Review. Drugs Aging 2017;34:891–905. https://doi.org/10.1007/s40266-017-0501-7.

[3] Kiefer DS, Chase JC, Love GD, Barrett BP. The Overlap of Dietary Supplement and Pharmaceutical Use in the MIDUS National Study. Evid Based Complement Alternat Med 2014;2014:823853. https://doi.org/10.1155/2014/823853.

[4] Brantley SJ, Argikar AA, Lin YS, Nagar S, Paine MF. Herb–Drug Interactions: Challenges and Opportunities for Improved Predictions. Drug Metab Dispos 2014;42:301–17. https://doi.org/10.1124/dmd.113.055236.

[5] Birer-Williams C, Gufford BT, Chou E, Alilio M, VanAlstine S, Morley RE, et al. A New Data Repository for Pharmacokinetic Natural Product-Drug Interactions: From Chemical Characterization to Clinical Studies. Drug Metab Dispos 2020;48:1104–12. https://doi.org/10.1124/dmd.120.000054.

[6] Nicholson DN, Greene CS. Constructing knowledge graphs and their biomedical applications. Computational and Structural Biotechnology Journal 2020;18:1414–28. https://doi.org/10.1016/j.csbj.2020.05.017.

[7] Paine MF, Shen DD, McCune JS. Recommended Approaches for Pharmacokinetic Natural Product-Drug Interaction Research: a NaPDI Center Commentary. Drug Metabolism and Disposition 2018;46:1041. https://doi.org/10.1124/dmd.117.079962.

[8] Wang T, Yang J, Xiao Y, Wang J, Wang Y, Zeng X, et al. DFinder: a novel end-to-end graph embedding-based method to identify drug–food interactions. Bioinformatics 2023;39:btac837. https://doi.org/10.1093/bioinformatics/btac837.

[9] Wang L, Tafjord O, Cohan A, Jain S, Skjonsberg S, Schoenick C, et al. SUPP.AI: finding evidence for supplement-drug interactions. Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations, Online: Association for Computational Linguistics; 2020, p. 362–71. https://doi.org/10.18653/v1/2020.acl-demos.41.

[10] Lin K, Friedman C, Finkelstein J. An automated system for retrieving herb-drug interaction related articles from MEDLINE. AMIA Jt Summits Transl Sci Proc 2016;2016:140–9.

[11] Trinh K, Pham D, Le L. Semantic Relation Extraction for Herb-Drug Interactions from the Biomedical Literature Using an Unsupervised Learning Approach. 2018 IEEE 18th International Conference on Bioinformatics and Bioengineering (BIBE), 2018, p. 334–7. https://doi.org/10.1109/BIBE.2018.00072.

[12] Schutte D, Vasilakes J, Bompelli A, Zhou Y, Fiszman M, Xu H, et al. Discovering novel drug-supplement interactions using SuppKG generated from the biomedical literature. Journal of Biomedical Informatics 2022;131:104120. https://doi.org/10.1016/j.jbi.2022.104120.

[13] Zhu X, Gu Y, Xiao Z. HerbKG: Constructing a Herbal-Molecular Medicine Knowledge Graph Using a Two-Stage Framework Based on Deep Transfer Learning. Front Genet 2022;13:799349. https://doi.org/10.3389/fgene.2022.799349.

[14] Rahman MM, Vadrev SM, Magana-Mora A, Levman J, Soufan O. A novel graph mining approach to predict and evaluate food-drug interactions. Sci Rep 2022;12:1061. https://doi.org/10.1038/s41598-022-05132-y.

[15] He X, Zhang R, Rizvi R, Vasilakes J, Yang X, Guo Y, et al. Prototyping an Interactive Visualization of Dietary Supplement Knowledge Graph. Proceedings (IEEE Int Conf Bioinformatics Biomed) 2018;2018:1649–52. https://doi.org/10.1109/BIBM.2018.8621340.

[16] He X, Zhang R, Rizvi R, Vasilakes J, Yang X, Guo Y, et al. ALOHA: developing an interactive graph-based visualization for dietary supplement knowledge graph through user-centered design. BMC Med Inform Decis Mak 2019;19:150. https://doi.org/10.1186/s12911-019-0857-1.

[17] Zhang R, Hristovski D, Schutte D, Kastrin A, Fiszman M, Kilicoglu H. Drug repurposing for COVID-19 via knowledge graph completion. Journal of Biomedical Informatics 2021;115:103696. https://doi.org/10.1016/j.jbi.2021.103696.

[18] Tripodi IJ, Callahan TJ, Westfall JT, Meitzer NS, Dowell RD, Hunter LE. Applying knowledge-driven mechanistic inference to toxicogenomics. Toxicology in Vitro 2020;66:104877. https://doi.org/10.1016/j.tiv.2020.104877.

[19] Stockley's Herbal Medicines Interactions. MedicinesComplete n.d. https://about.medicinescomplete.com/publication/stockleys-herbal-medicines-interactions-2/ (accessed March 31, 2023).

[20] Bordea G, Nikiema J, Griffier R, Hamon T, Mougin F. FIDEO: food interactions with drugs evidence ontology. 11th International Conference on Biomedical Ontologies, 2020.

[21] Callahan TJ, Tripodi IJ, Hunter LE, Baumgartner WA. A Framework for Automated Construction of Heterogeneous Large-Scale Biomedical Knowledge Graphs. bioRxiv 2020:2020.04.30.071407. https://doi.org/10.1101/2020.04.30.071407.

[22] Callahan T. PheKnowLator 2019. https://doi.org/10.5281/zenodo.3401437.

[23] Taneja SB, Callahan TJ, Brochhausen M, Paine MF, Kane-Gill SL, Boyce RD. Designing potential extensions from G-SRS to ChEBI to identify natural product-drug interactions, 2021. https://doi.org/10.5281/zenodo.5736386.

[24] Kilicoglu H, Rosemblat G, Fiszman M, Shin D. Broad-coverage biomedical relation extraction with SemRep. BMC Bioinformatics 2020;21:188. https://doi.org/10.1186/s12859-020-3517-7.

[25] Gyori BM, Bachman JA, Subramanian K, Muhlich JL, Galescu L, Sorger PK. From word models to executable models of signaling networks using automated assembly. Molecular Systems Biology 2017;13:954. https://doi.org/10.15252/msb.20177651.

[26] Valenzuela-Escárcega MA, Babur Ö, Hahn-Powell G, Bell D, Hicks T, Noriega-Atala E, et al. Large-scale automated machine reading discovers new cancer-driving mechanisms. Database (Oxford) 2018;2018. https://doi.org/10.1093/database/bay098.

[27] Taneja SB, Callahan TJ, Paine MF, Kane-Gill SL, Kilicoglu H, Joachimiak MP, et al. Developing a Knowledge Graph for Pharmacokinetic Natural Product-Drug Interactions. Journal of Biomedical Informatics 2023;140:104341. https://doi.org/10.1016/j.jbi.2023.104341.

[28] Taneja SB. NP-KG: Knowledge Graph Framework to Generate Hypotheses for Natural Product-Drug Interactions (v1.0.1) 2022. https://doi.org/10.5281/zenodo.7011488.