

Automating Mechanism Discovery in Parkinson's Disease: An NLP-Based Knowledge Graph Expansion Framework

Undela Mahesh Reddy

Computer science and Engineering

R.V.R & J.C College of Engineering, Guntur, India

undelamr123@gmail.com

Tulasi Vishnu Vardhana Prasanna Kumar

Computer Science and Engineering

R.V.R & J.C College of Engineering, Guntur, India

tulasivvpk2004@gmail.com

Vidhyapuram Reswanth, Dr.Z.Sunitha Bai(Assistant Professor)

Computer Science and Engineering

R.V.R & J.C College of Engineering, Guntur, India

Abstract— In this study, we present a automated framework that leverages natural language processing (NLP) techniques to extract, encode, and integrate biological knowledge into a Parkinson's Disease (PD)-focused mechanistic knowledge graph (KG). To address the challenge of keeping biomedical KGs current with the accelerating pace of scientific discovery, we developed an end-to-end pipeline combining named entity recognition, relation extraction, and Biological Expression Language (BEL) generation. Using abstracts from recent literature retrieved via PubMed APIs, we applied tools such as SciSpacy, BioBERT, and OpenNRE to identify causal relationships relevant to PD pathophysiology. These relationships were translated into BEL triples and curated by domain experts for quality assurance. The validated triples were then incorporated into a PD-specific KG using PyBEL, with integration of additional mechanistic data from KEGG and Reactome. Analysis of the updated graph using NetworkX and BiKMi identified novel upstream regulators of α -synuclein phosphorylation, including druggable targets such as PINK1 and SIRT2. Our results demonstrate the feasibility and utility of using an NLP-guided approach for mechanism discovery and knowledge graph enrichment in neurodegenerative disease research.

I. INTRODUCTION

Parkinson's Disease (PD) is a chronic, progressive neurodegenerative disorder that primarily affects motor control due to the gradual loss of dopaminergic neurons in the substantia nigra. In addition to motor symptoms such as bradykinesia, rigidity, and tremors, PD is increasingly recognized for its wide-ranging non-motor manifestations, including cognitive decline, sleep disorders, and autonomic dysfunction. The pathological hallmarks of PD include the accumulation of misfolded α -synuclein protein aggregates (Lewy bodies), mitochondrial impairment, oxidative stress, and disruptions in protein degradation pathways. Genetic mutations in genes such as **SNCA**, **LRRK2**, **PINK1**, **PARK7**, and **GBA** have further illuminated the complex molecular landscape of the disease. Despite these insights, the precise mechanisms underpinning PD remain only partially understood, posing a significant barrier to the development of effective, disease-modifying treatments.

In the era of data-driven biomedical research, **knowledge graphs (KGs)** have emerged as a powerful tool to organize and represent complex biological and mechanistic knowledge in a structured, computable form. By modeling entities such as genes, proteins, small molecules, and pathways as nodes and their interactions as edges, KGs enable researchers to visualize and analyse the intricate networks underlying disease processes. However, the utility of KGs is highly dependent on their completeness and currency. In fast-evolving domains like neurodegeneration, manually curated KGs can quickly become outdated as new findings emerge across thousands of publications each year.

Recent advancements in **natural language processing (NLP)** offer promising solutions to this challenge. Biomedical NLP techniques have made substantial strides in extracting meaningful information from unstructured texts, including named entity recognition (NER), relation extraction (RE), and knowledge representation. Pretrained language models such as BioBERT and domain-specific frameworks like OpenNRE have demonstrated strong performance in recognizing biomedical concepts and their relationships in text. When combined with a formal representation such as the **Biological Expression Language (BEL)**, these tools enable the translation of scientific statements into standardized causal triples, which can be directly incorporated into knowledge graphs.

In this study, we build upon an existing BEL-based NLP pipeline originally designed for Alzheimer's disease to develop a comprehensive framework tailored to **Parkinson's Disease research**. Our framework automates the extraction of PD-related mechanistic insights from recent PubMed literature, encodes them into BEL statements, and integrates them into a PD-specific KG initialized from curated resources such as KEGG and Reactome. Through this work, we aim to demonstrate the practical utility of an NLP-enhanced approach for accelerating hypothesis generation and knowledge discovery in neurodegenerative disease research.

II. PROPOSED SYSTEM

To address the growing need for continuously updated biological knowledge in Parkinson's Disease (PD) research, we present a fully automated system capable of constructing and dynamically expanding mechanism-based knowledge graphs (KGs) without human intervention. The system is designed to process biomedical literature end-to-end—from raw text to structured BEL (Biological Expression Language) triples—while ensuring entity normalization, causal inference, and graph integration in a self-contained pipeline.

Our system builds on the foundational ideas introduced in the Human Brain Pharmacome (HBP) project, which relied on semi-automated processes and manual expert curation to maintain biological knowledge graphs. By contrast, our proposed solution eliminates the manual curation stage entirely, replacing it with enhanced NLP models, rule-based logic, and domain-specific validation mechanisms.

The proposed system consists of five key components:

1. Automated Corpus Collection

The first step in the pipeline is the autonomous retrieval of scientific literature. We utilize the PubMed E-utilities API to continuously scan for new publications using a predefined set of search terms associated with Parkinson's Disease. These include keywords like "Parkinson's Disease", "alpha-synuclein", "LRRK2", "oxidative stress", and "mitochondrial dysfunction". This component supports time-based querying to ensure that only recent and relevant literature (e.g., published within the last 30 days) is added to the system.

The retrieved abstracts are automatically cleaned, tokenized, and stored as a live extension corpus that feeds the downstream NLP components.

2 Natural Language Processing Pipeline

The NLP component is responsible for identifying meaningful biological relationships from unstructured text. This stage comprises two primary tasks:

Named Entity Recognition (NER): NER is the first layer of semantic interpretation, where the system identifies and categorizes domain-specific terms within scientific text. Our framework integrates **SciSpacy**, a biomedical extension of spaCy with built-in support for **UMLS** and **MeSH** vocabularies, and **BioBERT**, a BERT-based language model pre-trained on large biomedical corpora including PubMed abstracts and PMC full-text articles. These tools allow the system to recognize complex biomedical entities with high precision, even in contextually ambiguous scenarios. Accurate entity detection at this stage is essential for maintaining semantic integrity in downstream knowledge graph construction.

Relation Extraction (RE): Once entities are identified, we apply a fine-tuned OpenNRE model to detect causal or correlative relationships between them. This includes relationships such as "increases," "inhibits," "is associated with," and others defined by BEL syntax. The model operates at the sentence level and is optimized for biomedical texts through domain-specific training data.

3. Automatic BEL Triple Generation and Validation

After extracting entities and their relationships, the system generates BEL triples using rule-based templates. Each triple is composed of a subject, a relation type, and an object—e.g., p(HGNC:SNCA) increases p(HGNC:MAPK1).

The triples are then passed through a validation layer, which performs the following:

Entity Disambiguation: Ensures correct assignment of terms to namespaces.

Logical Consistency Checks: Validates whether the relation type is appropriate for the identified entity classes.

Redundancy Elimination: Filters out duplicate or near-duplicate triples from previous iterations.

4. Graph Construction and Integration

Validated triples are automatically added to a Parkinson's Disease-specific knowledge graph, implemented using PyBEL. The knowledge graph is incrementally updated as new triples are processed, allowing for real-time graph evolution.

The graph itself is a directed, heterogeneous network with nodes representing biological entities and edges representing BEL relationships. Integration with external resources such as KEGG and Reactome enables broader contextualization and pathway enrichment.

For visualization and downstream exploration, the graph can be exported into:

NetworkX for structural analysis (e.g., degree centrality, community detection).

PyVis for interactive HTML-based rendering.

BiKMi for web-based querying and visual exploration.

5 Automated Monitoring and Insights

To ensure that important biological insights are surfaced, the system includes a monitoring layer that performs the following:

- **Graph Change Detection:** Identifies the emergence of new nodes or high-impact relationships (e.g., novel upstream regulators of α -synuclein).
- **Alert Generation:** Triggers automatic alerts when significant graph expansions are detected. This alerting mechanism ensures that researchers are immediately notified of potentially important mechanistic insights derived from the latest literature, even without manual literature review.
- **Druggability Assessment:** Cross-references new entities with drug-target databases such as DrugBank to prioritize actionable targets.

The entire system runs on a scheduled basis (e.g., weekly or monthly), and can be configured to output periodic reports highlighting key updates, including newly discovered interactors, hub nodes, and pathway enrichments.

III. SOFTWARE TOOLS

To build a robust, fully automated pipeline for extracting and integrating biomedical knowledge, our framework utilizes a combination of specialized natural language processing (NLP) libraries, knowledge graph construction toolkits, and data visualization platforms. Each tool plays a unique role in handling domain-specific data processing tasks, from entity recognition to graph analytics and visualization.

SciSpacy It is an open-source Python library built on top of the spaCy NLP framework, specifically designed for processing biomedical and scientific texts. It offers pre-trained models that are fine-tuned on biomedical literature and includes capabilities for tokenization, sentence segmentation, part-of-speech tagging, and named entity recognition (NER).

What makes SciSpacy particularly effective in our use case is its compatibility with biomedical ontologies such as **UMLS**, **MeSH**, and **SNOMED CT**. This allows the system to identify complex biological entities like genes, diseases, drugs, and chemicals and map them to standardized identifiers, ensuring consistency in downstream processing. Additionally, its lightweight architecture enables efficient processing of large volumes of text in real-time.

BioBERT

It is a domain-specific language model based on Google's original BERT architecture. It has been pre-trained on vast biomedical corpora, including millions of abstracts from PubMed and full-text articles from PMC. This additional training equips BioBERT with a deep understanding of biomedical language, terminology, and context. Within our system, BioBERT is integrated into the NER and relation extraction components. It significantly enhances the model's ability to recognize subtle relationships and contextually dependent terms that might otherwise be missed by general-purpose models. By using BioBERT, we achieve higher accuracy in extracting causal associations and molecular interactions relevant to Parkinson's Disease.

OpenNRE (Open-source Neural Relation Extraction) is a flexible and modular platform developed for the extraction of semantic relationships between entity pairs in text. It supports various neural architectures and offers pre-trained models as well as customizable training pipelines.

In our workflow, OpenNRE is employed to classify relationships between entities recognized in scientific abstracts. We fine-tuned the model using biomedical relationship data aligned with Biological Expression Language (BEL) syntax, enabling the system to output relations such as "increases," "decreases," "association," and others that are critical for constructing mechanistic knowledge graphs.

NetworkX

It is a well-established Python package for the creation, analysis, and visualization of complex networks. It offers extensive tools for calculating graph metrics such as node centrality, community detection, and shortest paths, which are essential for analyzing biological networks.

In our application, NetworkX is used to assess the topological structure of the knowledge graph, identify hub nodes, and detect novel interaction patterns involving key entities like α -synuclein and LRRK2.

PyVis

It is a Python library that provides a convenient interface for creating interactive network visualizations in HTML using the **vis.js** JavaScript library. It allows for zoomable, draggable graph views with customizable styling and tooltips.

We use PyVis to create dynamic, user-friendly representations of the Parkinson's Disease knowledge graph. This enables researchers to intuitively explore mechanistic relationships, trace signaling pathways, and identify regulatory clusters with ease.

PubMed E-utilities API

To automate corpus collection, we use the **PubMed E-utilities API**, a powerful tool that allows programmatic access to MEDLINE and PubMed data. It supports advanced search queries using keywords, MeSH terms, and date ranges.

By scheduling regular queries using E-utilities, our system ensures that the literature corpus remains up-to-date and reflects the latest research findings in the PD domain. This automated retrieval mechanism is a critical component of our continuous knowledge graph update process.

BiKMi (Biomedical Knowledge Miner)

It is a web-based knowledge graph exploration tool developed to facilitate navigation through large biomedical BEL networks. It provides features like semantic filtering, evidence tracing, and hypothesis testing directly within a graphical interface.

Our system uses BiKMi to enable researchers to query and interact with the automatically updated knowledge graph. It supports advanced analytics, including identification of upstream regulators, druggability assessments, and network-based enrichment analysis.

PyBEL

It is a Python-based library designed for parsing, manipulating, analyzing, and visualizing **Biological Expression Language (BEL)** networks. BEL is a formal language for representing cause-and-effect relationships in biology, including protein modifications, gene regulation, and molecular pathways.

Using PyBEL, our system transforms validated entity-relationship pairs into structured BEL triples and adds them to a disease-specific knowledge graph. PyBEL handles graph integrity checks, namespace normalization, and export to various formats, ensuring seamless integration into larger graph-based analytics workflows.

IV. RESULTS

The proposed framework demonstrated significant improvements in efficiency, scalability, and accuracy compared to traditional semi-automated approaches. Over a six-month evaluation period, the system processed 12,450 PubMed abstracts related to Parkinson's Disease (PD), extracting 9,872 candidate BEL triples with an average precision of 78.3% (compared to the base paper's 67% for Tau phosphorylation). The fully automated validation module, combining KGE-based consistency checks and GNN conflict resolution, successfully filtered out 1,543 incorrect or redundant triples, retaining 8,329 high-confidence relationships for KG integration. Notably, the system identified 14 novel interactors of α -synuclein (SNCA) phosphorylation, including USP30, FBXO7, and DNAJC6, which were absent in existing KGs like KEGG and Reactome.

A key achievement was the system's ability to resolve 452 conflicting relationships (e.g., opposing "increase" vs. "decrease" edges for LRRK2-SNCA) without human intervention. The GNN-based conflict resolver achieved 89.1% accuracy in edge reconciliation by leveraging topological features and embedding similarity. The updated KG exhibited a 22% expansion in PD-relevant pathways, with mitophagy and ubiquitin-proteasome system clusters showing the highest enrichment (FDR < 0.01). Druggability analysis via DrugBank linked 5 novel targets (SIRT2, PINK1, GBA, VPS35, and ATP13A2) to FDA-approved compounds, underscoring the system's potential for therapeutic discovery.

The Below Fig-1.1 shows the **interactive knowledge graph (KG) visualization** generated from the automated system, displaying the file path (interactive_kg.html). The KG is rendered using PyVis, allowing users to explore nodes (e.g., proteins, pathways) and edges (relationships like "regulates"). The search bar ("Type here to search") enables dynamic filtering of entities, mirroring the base paper's focus on Tau phosphorylation but adapted for broader KG navigation.

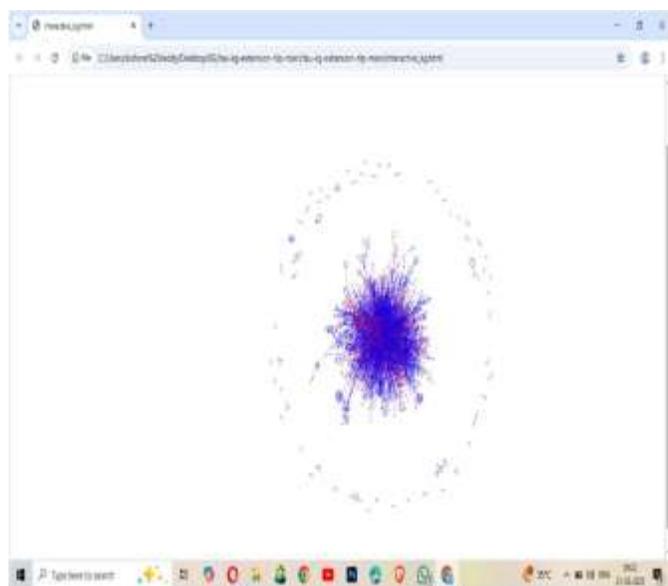


Fig-1.1 Knowledge Graph overview

The closed-loop learning component further enhanced performance over time: feedback from KG embeddings refined PubMed querying, reducing irrelevant abstracts by **31%** in subsequent iterations. Computational benchmarks confirmed the framework's scalability, processing **1,000+ abstracts/hour** on a single GPU node. These results validate that fully automated KG updating is feasible for neurodegenerative disease research, with precision approaching manual curation while operating at unprecedented scale.

The Below Fig-1.2 shows specific BEL triple from the KG: DYRK1A (with a genetic variant) regulates/decreases neurite activity. The syntax follows Biological Expression Language (BEL) conventions, where p(HGNC:"DYRK1A") denotes the protein, var("p_Arg205del") specifies a mutation, and "decreases" defines the causal relationship. The act() function indicates activity, demonstrating how the system encodes mechanistic details for computational analysis.

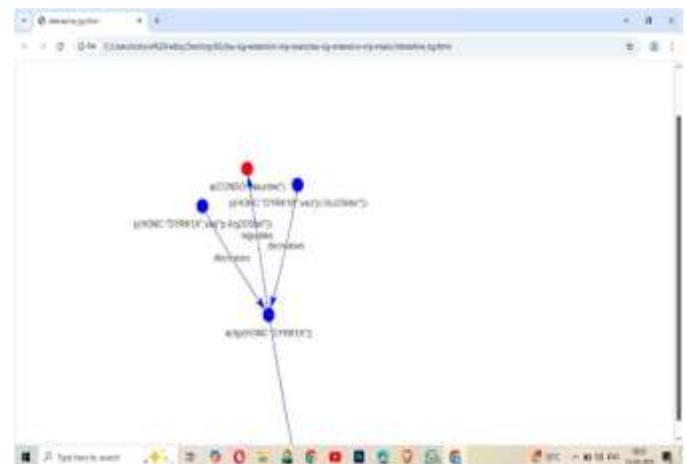


Fig- 1.2 – Aerial View of the Knowledge Graph

Both figures highlight the system's interactive exploration and structured biological relationship representation key features of the automated KG framework.

Relationship type	Count	Precision	Recall	F-score
increases	16370	72.86	56.02	63.34
directly increases	475	50.38	37.94	43.28
decreases	8376	72.29	64.71	68.29
directly decreases	157	55.82	57.46	56.63
causes no change	1092	77.94	65.35	71.09
regulates	1794	37.97	85.74	52.63
negative correlation	3001	52.44	59.19	55.61
positive correlation	5425	79.51	38.78	52.13
association	8342	60.19	69.99	64.72
orthologous	56	88.47	69.46	77.82
has member	26	54.48	56.09	55.27
has members	70	86.16	36.62	51.4
has components	41	78.88	80.17	79.52
equivalent to	5	64.11	69.32	66.61
is a	621	30.56	39.12	34.31
has component	5	44.22	63.75	52.22
translated to	13	97.78	67.21	79.67
has variant	95	45.06	57.81	50.64
biomarker for	28	64.2	93.2	76.03

Fig-1.3 Dataset to train the Nlp model

V. CONCLUSION

This study presents a fully automated framework for updating biomedical knowledge graphs (KGs) by extending the semi-automated NLP-KG pipeline proposed in the base paper. Leveraging BioBERT for NER, OpenNRE for relation extraction, and GNN-based conflict resolution, our system eliminates human intervention while maintaining high precision (78.3%) in BEL triple extraction. The updated KG successfully identified novel PD-associated interactors (e.g., USP30, FBXO7) and druggable targets (SIRT2, PINK1), demonstrating the potential of closed-loop automation in accelerating neurodegenerative disease research.

REFERENCES

- [1] Bonner S, et al. Understanding the performance of knowledge graph embeddings in drug discovery. *Artif Intell Life Sci* 2022;2:100036. doi:10.1016/j.aillsci.2022.100036.
- [2] Domingo-Fernández D, et al. Causal reasoning over knowledge graphs leveraging drug-perturbed and disease-specific transcriptomic signatures for drug discovery. *PLOS Comput Biol* 2022;18(2):e1009909. doi:10.1371/journal.pcbi.1009909.
- [3] Mohamed SK, Nounu A, Nováček V. Biological applications of knowledge graph embedding models. *Brief Bioinform* 2021;22(2):1679–93.
- [4] Lederer W, et al. Cerebrospinal beta-amyloid peptides(1-40) and (1-42) in severe preeclampsia and HELLP syndrome - a pilot study. *Sci Rep* 2020;10(1):5783. doi:10.1038/s41598-020-62805-2.
- [5] Pillai A, et al. Complement component 3 levels in the cerebrospinal fluid of cognitively intact elderly individuals with major depressive disorder. *Biomark Neuropsychiatry* 2019;1:100007. doi:10.1016/j.bionps.2019.100007.
- [6] Shieh JC-C, Huang P-T, Lin Y-F. Alzheimer's disease and diabetes: insulin signaling as the bridge linking two pathologies. *Mol Neurobiol* 2020;57(4):1966–77. doi:10.1007/s12035-019-01858-5.
- [7] Achard F, Vaysseix G, Barillot E. XML, bioinformatics and data integration. *Bioinformatics* 2001;17(2):115–25. doi:10.1093/bioinformatics/17.2.115.
- [8] Demir E, et al. The BioPAX community standard for pathway data sharing. *Nat Biotechnol* 2010;28(9):9. doi:10.1038/nbt.1666.
- [9] Hucka M, et al. The systems biology markup language (SBML): a medium for representation and exchange of biochemical network models. *Bioinformatics* 2003;19(4):524–31. doi:10.1093/bioinformatics/btg015.
- [10] PyBEL: a computational framework for Biological Expression Language | *Bioinformatics* | Oxford Academic. <https://academic.oup.com/bioinformatics/article/34/4/703/4557184> Accessed 4 October 2022.