# KGWRAPS: A Knowledge Graph-Powered Research Assistant for Polymer Science

## 1  Problem Statement

Polymer science is a subfield of materials science concerned with polymers, substances or materials consisting of very large molecules that are in turn composed of several repeating subunits. It is a constantly evolving field, dedicated to the development of materials with improved performance, sustainability, and novel functionalities. Some advanced aspects of the field include the applications of smart polymers, polymers that are being widely used to revolutionize the development of smart devices, sensors, and actuators due to their ability to respond with a detectable reaction to environmental stimuli [7] and biodegradable polymers, polymers that degrade into natural byproducts after accomplishing their intended function [28].

Despite the importance of the field, however, various issues plague it and significantly limit the efficiency and robustness of various processes. For example, polymer science currently relies heavily on the traditional methods of manually controlled experimentation. As in other fields, this can lead to potential bias and human error in the design and analysis of the experiment. Another issue stemming directly from the excessive level of human involvement in experimentation and related closely to the field of polymer science is the reproducibility crisis. Scientists are finding it increasingly difficult to replicate the results of studies conducted by others. For example, as in the case of LK-99 (a supposed room-temperature superconductor), attempts at replicating the material with its superconductor-like properties failed, rendering uncertain whether the material was ever a superconductor to begin with and a mistake had simply been made [16]. Similar issues arise in the field of polymer chemistry as well. For example, many polymers possess a stochastic nature, comprising an ensemble of similar molecules rather than being well-defined single structures [3] and varying characteristics due to the polymerization method used for synthesis [27]. Thus, challenges to replication are a frequent issue plaguing the field. Additionally, with the reliance on manual experimentation arises the need to investigate a massive search space to gather the needed data. Because of the wealth of information from which this data must be drawn, this process can be very difficult and time-consuming.

These problems can be alleviated through the use of modern technology to automate several of the constituent processes of polymer science experimentation. The automation of these processes would significantly raise their efficiency by serving as a more optimal alternative to the traditional methods of rote experimentation. For example, the use of machine learning (ML) techniques can maximize experimental precision, scalability, and objectivity – factors that are difficult to maximize in an environment controlled solely by humans. In addition, such autonomous, computer-aided research techniques would assist in diminishing the impact of the reproducibility crisis, as they would allow for more control and consistency in the design of experiments and more detailed analysis of resulting data. Knowledge engineering, in particular, serves as a potent solution to the issue of a massive search space. Knowledge graphs (KGs), a key component of the field, offer a structured representation of data that allows convenient access to relevant information. KGs are already being used effectively in the realm of polymer science to make advancements. For example, Google DeepMind's Graph Networks for Materials Exploration (GNoME) is a graph neural network model being used to significantly increase the speed and efficiency of discovery of new materials by predicting their stability [11].

Thus, we propose the creation of Knowledge Graph-Powered Research Assistant for Polymer Science (KGWRAPS), a research assistant based on KGs to aid users in the development of experiments related to polymer science. The two key goals of the research assistant are to aid in experimentation related to the reproduction of already existing polymers and the discovery of new polymers. Through these objectives, it will also achieve accurate prediction of material performance and characterization based on the effect of factors such as polymer chemistry, state variables, and other thermodynamic properties. With the added benefits of autonomous, artificially intelligent systems, we anticipate increased efficiency and quality in the production of polymers assisted by this research assistant.

## 2 Background and Relevance

The implementation of this research assistant relies on knowledge drawn from a variety of fields: polymer science (Section 2.1), knowledge engineering (Section 2.2), and conversational artificial intelligence (AI) (Section 2.3).

### 2.1 Polymer Science

As the domain in which this research is to be conducted, polymer science plays a crucial role in its progression. The integration of polymer science with computer-aided technologies has done much to advance the state of the art of the field. ML, in particular, has been used extensively for optimization of techniques and processes in polymer science, as well as to allow researchers to uncover patterns and correlations through the processing of large amounts of data. Bayesian optimization, an application of ML, for example, has been used to optimize chemical reactions with the goal of achieving functionalization of polymers [13]. With polymer processing via solution phase behavior being a significant area of study within the field of polymer science and for the Air Force Research Laboratory's (AFRL) applications specifically, this work demonstrates the extent to which ML can be utilized in the modification of polymers and their properties.

In the context of materials science research, computer-aided technologies play a pivotal role in challenging the status quo. AFRL's ARES, the Autonomous Research System, for example, allows the transformation of automated experiments and carbon nanotube synthesis reactors into autonomous systems capable of directing and conducting their own research using AI. The purpose of this software is to make research faster and more effective by minimizing time spent on traditional research processes [21]. Additionally, it and other autonomous experimentation tools lower barriers to entry into scientific fields by reducing the level of resources that need to be invested in research [2].

### 2.2 Knowledge Engineering

Knowledge engineering and specifically KGs offer a structured representation of knowledge that organizes information in a way that is easily understandable by both humans and machines. A KG consists of entities, which are objects or concepts, and the relationships between them. These entities and relationships are typically represented as nodes and edges in a graph, respectively [19]. Figure 1 represents an ontology illustrating the modeling approach for processing polymer nanocomposites in the NanoMine KG. It demonstrates that each sample of polymer nanocomposite is generated by a process that has several steps. In addition to various attributes and equipment used, each step has inputs and outputs; and the use of one step's output as the next step's input provides a level of ordering to the overall procedure [8]. KGs are valuable for organizing and connecting information in a semantically meaningful way. They enable more effective data retrieval, reasoning, and inference. KGs find applications in various domains, such as natural language processing, AI, and data integration, helping to create a more comprehensive and interconnected understanding of information.

KGs have many well-known applications such as Google's Knowledge Graph and Wikidata. Related to materials science specifically, Google DeepMind's GNoME has predicted 2.2 million new crystals, equivalent to approximately 800 years' worth of knowledge. Of these materials, 380,000 are considered sufficiently stable, rendering experimental synthesis plausible [11]. NanoMine, discussed previously, is another KG in the realm of polymer science that integrates data from over 1,700 polymer nanocomposite experiments to assist in exploration of the effects of polymer nanocomposite design on their properties [8]. MatKG is an attempt at establishing a field-wide ontology on materials science. Using transformer-based large language models (LLMs), this graph database is composed of 2 million unique relationships among over 80,000 unique entities autonomously extracted from over 4 million papers on the topic of materials [33].

### 2.3 Conversational Artificial Intelligence

Conversational AI represents a groundbreaking advancement in technology that enables machines to engage in natural language conversations with users. This innovative field of AI leverages techniques such as natural language processing, ML, and deep learning to understand and respond to human inputs. The benefits of conversational AI are manifold, ranging from improved customer support and enhanced user experience to increased efficiency in various industries. By facilitating seamless interactions between humans and machines, conversational AI streamlines processes, provides quick and accurate information, and offers a personalized touch in digital interactions. It plays a pivotal role in automating tasks, reducing workload, and fostering a more intuitive and accessible digital environment for users across diverse applications, from virtual assistants and chatbots to customer service platforms and smart devices.

With regard to the realm of polymer science, conversational AI has the potential to make research easier and more accessible. It serves to enable the design of experiments and analysis of data in a way that is convenient for humans: natural language. One study conducted a hackathon event focused on understanding the applicability of LLMs to materials science and chemistry, concluding that the completed projects would have taken several months without the use of these LLMs [17].

## 3 Proposed Methodology

This project will be executed along five phases, outlined as follows. An overview of the components of the project is graphically depicted in Figure 2.

### 3.1 Data Discovery and Collection

We will first conduct a literature survey and dataset search. At first, we will focus on already-known information about the effects of factors, such as solvent quality, experimental conditions, polymer chemistry, and kinetics on material performance and characterization. We anticipate that we will obtain the data from online databases and datasets, but will need to be aware of trades-offs for available sources [5].

**Existing Sources.** Some prominent data sources include the Materials Project, which allows open access to information on known and predicted materials along with analysis tools to aid in the design of novel materials [20]; the Materials Genome Initiative (MGI) data repository, which focuses on material measurement and data delivery through the development, solving, and quantification of materials models and tools for materials theory and modeling [10]; and the PoLyInfo database, which provides various data such as properties, chemical structures, IUPAC names, processing methods of measured samples, measurement conditions, used monomers, and polymerization methods required for polymeric material design [26]. This list is simply a small subset of the numerous sources of materials information from which our research can draw.

### 3.2 Domain Knowledge Elicitation

Following the collection of the needed data, we will seek guidance from domain experts in the field of polymer science, particularly polymer scientists and engineers at AFRL. The purpose of this step is to begin encoding domain (polymer science) expertise into the KG. For KGWRAPS to support research, it must have its own epistemological model of the domain to connect ontologically to the relevant data. The faculty advisor has significant experience in these types of knowledge elicitation tasks. Furthermore, we will make use of the Toolbox Dialogue Initiative [32] to follow best-practices in bridging transdisciplinary gaps in vocabulary and understanding, which will drastically improve ontological outcomes.

### 3.3 Knowledge Graph Development

To create a maximally adaptable and reusable KG, we will use the Modular Ontology Modeling (MOMo) Methodology [31]. This methodology is an established best-practice, which has been used to create several

high-profile KGs [18, 14]. Sections 3.1 and 3.2 overlap with the initial steps – it is important to distinguish the slight differences in this case, due to the number of available resources (both data and expertise) available for this proposed work. The schema will outline at a high level the relationships among the data and its basic structure.

**Mapping from Existing Sources.** To map from the existing semantic resources to our modularly-developed schema, as outlined in Section 2.2, we will utilize [9], a formally-specified RDF [15] mapping language. R2RML's flexibility enables the automated creation of tailored mappings from relational schema to RDF that align seamlessly with the specific requirements of a KG. This adaptability is essential, as our knowledge representation demands nuanced relationships and structures.

**Tool Network & Sensors.** This aspect of the KG pertains to the integration of experimental and environmental conditions. These might be temperature, humidity, or air pressure, or other properties; this network of tools and sensors will be used to enhance our dataset with data collected in real time. Thus, not only will the schema need to incorporate the structure to be used for data collected from pre-existing databases and datasets, but it will also incorporate contextual data invaluable for replicability or otherwise understanding experimental results. Due to the real time nature of this aspect, we will use OpenRefine (which can be deployed as a server [22]) to on-the-fly map raw data into a KG through additional R2RML mappings [9].

## 3.4 Research Support

KGWRAPS' entire goal is to support and accelerate research, experiment design, and campaign strategy. To accomplish this, we will initially support the following functionalities.

**Source of Truth.** A KG SoT will enable straightforward querying for relevant data. This can be implemented through various tools, namely SPARQL [30], or bespoke interfaces. Additionally, the KGWRAPS system will utilize the contextualized KG as its memory, enabling additional research support functionalities.

**KG Embeddings & Graph Analytics.** KG Embeddings (KGEs) are a method of vectorizing KG data [4]. For example, we would then be able to observe clusters of polymers based on specific properties or environmental conditions. By doing this, the research assistant would be able to identify what these materials have in common and even how to produce materials with similar or more pronounced properties. Finally, the "slices" of the KG can be extracted in various formats to enable well-known graph analytics and network science algorithms [24].

**Knowledge Gap Identification.** We intend for KGWRAPS to be capable of making predictions about polymer properties; however, we also want to identify potential gaps in our knowledge based on the contents of our KG SoT. We will build a hybrid approach to detecting gaps in available knowledge, starting with work in [29]. That is, we will use heuristics for searching for gaps in the SoT but also examine the latent space of the KGE over the SoT. For example, we will be able to make predictions about polymer properties based on proposed experimental conditions. These predictions will be especially helpful in the production of polymers, as they will limit investment of time and resources needed to find relationships between material properties and experimental conditions. In addition, by filling these knowledge gaps, the research assistant can use new knowledge in future predictions and the design of future experiments.

**Conversational Experimental Design.** This component builds on the previous two functionalities. After KGWRAPS has identified a knowledge gap in the SoT, KGWRAPS will be capable of assisting the research scientist in designing an appropriate experiment to fill that knowledge gap. We will adapt [12] from the medical conversation domain to the material science domain – this turns out to be quite useful because it encodes intents but also seeks to confirm conversation state quite frequently, which is just as important in medicine as it is in designing expensive experiments. This process will populate a subgraph in the SoT, structured according to our schema developed in Sections 3.2 and 3.3.

### 3.5 Evaluation

**KG SoT Assessment.** This is a combined evaluation of both the ontology, the KG, and its ability to incorporate domain science. The ontology and KG will be evaluated according to best practices, which are included in the development methodology [31].

**Holistic Usability Assessment.** This evaluation regards the entire system and its usability. KGWRAPS will be assessed to determine the extent to which users find it practical and effective in their research endeavors, and evaluated along two metrics: technical correctness and ease of use. With regard to the technical correctness, this will require the close observation of domain research scientists to validate identified knowledge gaps, as well as valid experimental design.

The usability assessment will involve gathering feedback from users who have interacted with the research assistant, following the guidelines in [1], for a variety of metrics, including user satisfaction with KGWRAPS assistance in polymer science research, user efficiency in designing experiments, and practical utility of the research assistant in real-world polymer science research scenarios.

## 4 Expected Results and Significance

The implementation of KGWRAPS is anticipated to yield significant advancements in the field of polymer science. By automating various processes of polymer science experimentation, such as experiment design and analysis, the research assistant aims to address existing challenges associated with manual experimentation. The use of ML techniques within the research assistant is expected to enhance experimental precision, scalability, and objectivity, mitigating potential biases and errors introduced by human involvement.

The research assistant's focus on the reproduction of existing polymers and the discovery of new polymers aligns with the broader goals of polymer science, which include developing materials with improved performance, sustainability, and novel functionalities. Through accurate prediction of material performance based on factors such as polymer chemistry and thermodynamic properties, the research assistant aims to contribute to the production of polymers with enhanced properties and characteristics. The integration of ML and other computer-aided technologies is rapidly expanding in the field and resulting in unprecedented advances. Our work, seeking to bridge the gap between polymer science and knowledge engineering, is expected to further accelerate the rate at which these advances are occurring.

The incorporation of conversational AI in the research assistant is anticipated to streamline interactions between users and the system. This aspect holds the potential to make polymer science research more accessible by allowing users to engage with the assistant in natural language. The extensive benefits of LLMs in fields like chemistry, particularly the autonomous design, planning, and performance of complex scientific experiments, have been demonstrated to much success [6]. Thus, with the added benefit of a structured repository of information in the form of a KG, our research assistant's ability to assist users in experiment design, knowledge gap identification, and graph analytics is expected to greatly enhance the efficiency of research endeavors in polymer science.

In addition to the significant advancements anticipated in the field of polymer science through the implementation of our research assistant, we expect several high-impact publications and conference presentations to stem from the success of this project. These publications and presentations will serve as a testament to the collaborative effort between our team and AFRL, showcasing the accelerated design of new polymer materials and highlighting the innovative integration of modern technologies such as KGs and conversational AI. To accelerate advancement in polymer research through community involvement, we intend for KGWRAPS to be an open-source tool, licensed under the GNU Lesser General Public License, and will provide a link to its repository in these publications. By disseminating our findings through esteemed journals and conferences, we aim to not only contribute to the scientific community's understanding of polymer science but also foster continued collaboration in this rapidly evolving field.

## 5 Outline of Anticipated Efforts Beyond Year 1 & Special Requirements

Beyond the first year of the project, the research will focus on adding conversational capabilities to the research assistant. This process will involve combining an LLM with Convology. Convology (CONVersational ontOLOGY) is a top-level ontology aiming to model the conversation scenario for supporting the development of conversational knowledge-based systems. It defines concepts describing various components integral to a conversation and can theoretically be used in a multitude of disciplines. For example, Puffbot is introduced as a multiturn goal-oriented conversational agent based on Convology that supports patients affected by asthma [12]. By using Convology, we simplify the process of integrating natural language into our research assistant while also rendering it more accessible and easier to use.

Additionally, in light of the comprehensive scope outlined in Section 3, it is anticipated that certain aspects of the work may necessitate more time for thorough completion beyond the initial year. For example, a true, final evaluation of the research assistant cannot be completed before the conversational capabilities have been implemented. Figure 3 depicts a tentative schedule for the overarching, significant steps of this project.

**Special Requirements.** The work outlined above requires close collaboration with domain experts, which will only be available by working onsite. Additionally, access to a physical laboratory will allow us to validate KGWRAPS' behavior.

## 6 Potential Impact on Air Force and Ohio Economy

We anticipate KGWRAPS to have several lasting effects on the Air Force and Ohio economy. AFRL has a vested interest in polymer science due to its critical role in developing novel aerospace materials, enabling advancements in soft robotics, and facilitating breakthroughs in nanoelectronics with the goal of enabling seamless integration at the interface between warfighters and the machines they operate [25]. Polymers serve as essential components in a wide array of Air Force applications, including aircraft construction, protective coatings, and electronic devices. However, the traditional methods of polymer development and experimentation are time-consuming and resource-intensive, often leading to lengthy development cycles and limited innovation. By accelerating the discovery of polymer materials through the integration of AI and KGs, our project aims to revolutionize materials development processes at AFRL. This research assistant will leverage AI techniques to analyze vast amounts of data, identify correlations, and predict material properties with unprecedented accuracy and efficiency. By automating experimentation and data analysis, KGWRAPS will streamline the research process, enabling researchers to focus on high-impact tasks such as innovation and design. Moreover, the insights gained from this research assistant will not only lead to the rapid identification of materials with desired properties but also guide future experimentation, optimizing resource allocation and enhancing research efficiency at AFRL.

Due to the various products in which they can be found, polymers have become commonplace, especially in Ohio's industries. As such, Ohio contains the largest polymer industry cluster in America, being the state to produce the most plastic and rubber products, and is recognized as a global leader in this aspect [23]. Our research, seeking to revolutionize polymer science through the autonomous, more efficient development of robust polymers and polymer experiments will play a pivotal role in ensuring that Ohio maintains its role as a global leader in polymer innovation and production. By accelerating the discovery and development of advanced polymer materials, our research will not only bolster the competitiveness of Ohio's polymer industry but also contribute to the state's overall economic growth. Furthermore, the integration of AI and KGs in polymer science research will attract talent and investment to Ohio, fostering a vibrant ecosystem of innovation and entrepreneurship. Through collaboration with local industry partners and academic institutions, our research will catalyze the translation of scientific discoveries into practical applications, creating new opportunities for job creation and economic development in Ohio. As a result, our efforts will not only benefit the United States Air Force by advancing materials research for a wide range of applications crucial to national defense, but also contribute to the prosperity of Ohio's economy and its position as a global leader in polymer science and technology.
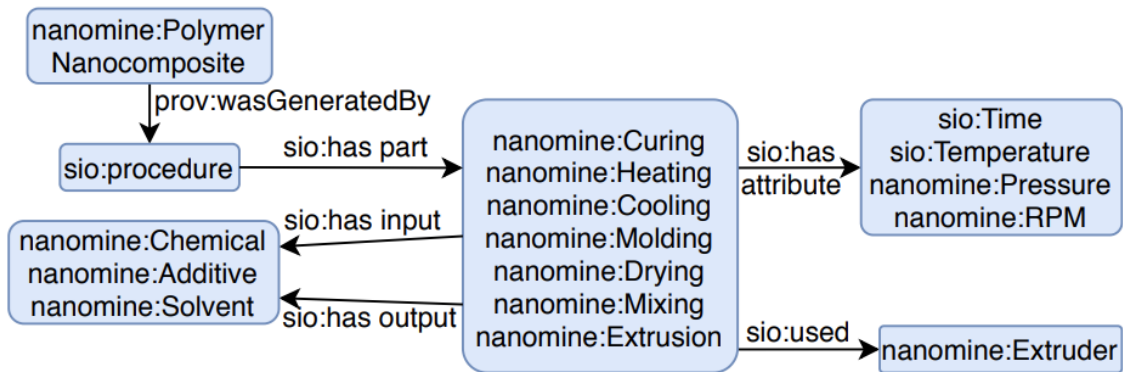
## A Figures



Figure 1: Excerpt of a knowledge graph from NanoMine: Ontological Representation of Modeling Processing of Polymer Nanocomposites [8]
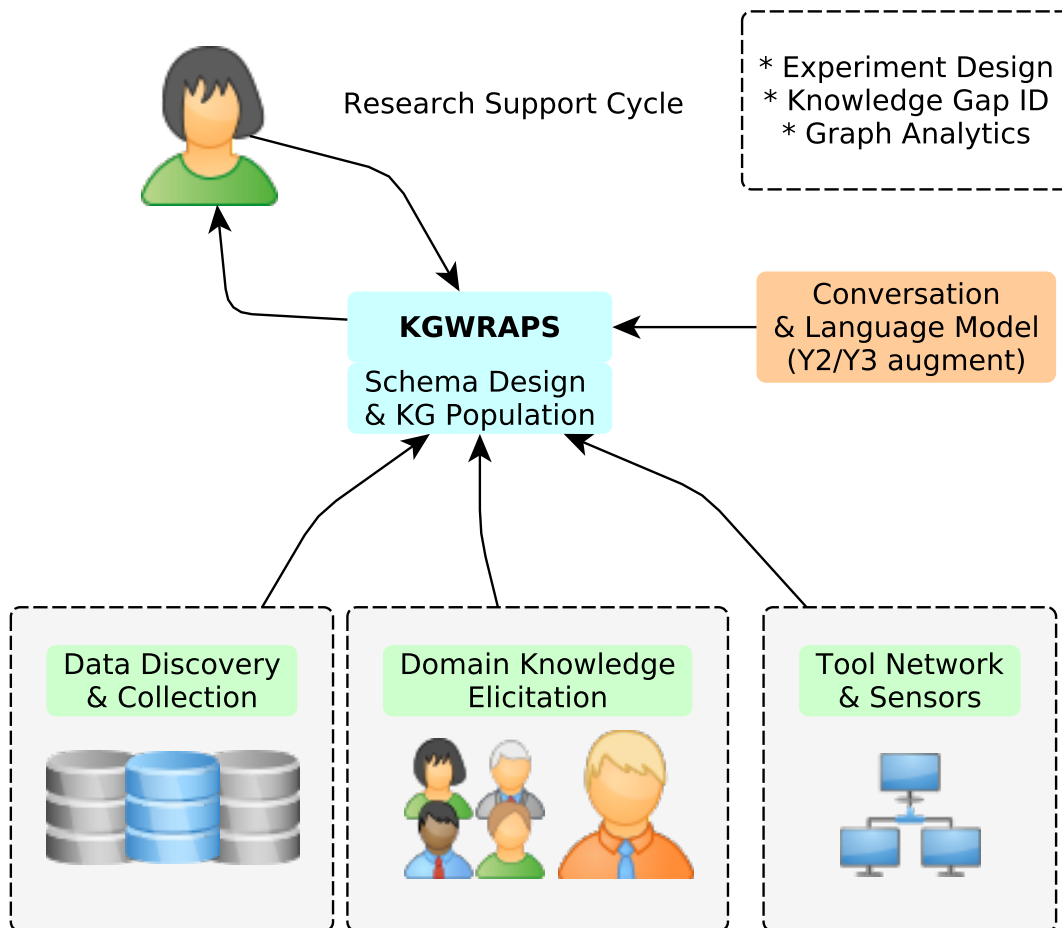


Figure 2: KGWRAPS: Visualization of methodology overview

| Task List | | | Schedule | | |
|---|---|---|---|---|---|
| Task # | Objectives | Tasks | Year 1 | Year 2 | Year 3 |
| 1 | Data Discovery and Collection | Integrate data from semantically enriched resources | May–Aug | | |
| 2 | | Integrate data from non-semantically enriched resources | May–Aug | | |
| 3 | Domain Knowledge Elicitation | Work with experts at AFRL to modify and organize data | Jun–Jan | | |
| 4 | Knowledge Graph Development | Design the schema, structuring and incorporating already collected data | Aug–Nov | | |
| 5 | | Modify the schema, structuring and incorporating data from the tool network and sensors | Nov–Dec | | |
| 6 | | Construct the KG based on the schema | Jan–Feb | | |
| 7 | Research Support | Implement graph analytics | | May–Jun | |
| 8 | | Implement knowledge gap identification | | Sept–Oct | |
| 9 | | Implement experiment design | | Jan–Feb | |
| 10 | Conversational AI Integration | Incorporate an LLM and Convology | | | May–Dec |
| 11 | Evaluation | Conduct a KG SoT assessment | | | Feb–Apr |
| 12 | | Conduct a holistic usability assessment | | | Feb–Apr |

Figure 3: A Gantt chart showing our anticipated timeline

## References

[1] A. S. f. P. Affairs. System Usability Scale (SUS). https://www.usability.gov/how-to-and-tools/methods/system-usability-scale.html, Sept. 2013. Publisher: Department of Health and Human Services.

[2] AI research robots key to 'democratizing and revolutionizing science,' world-class AFRL re. https://www.af.mil/News/Article-Display/Article/3562629/ai-research-robots-key-to-democratizing-and-revolutionizing-science-world-class, Oct. 2023.

[3] M. Aldeghi and C. W. Coley. A graph representation of molecular ensembles for polymer property prediction. *Chemical Science*, 13(35):10486–10498, 2022. arXiv:2205.08619 [cond-mat].

[4] M. Ali, M. Berrendorf, C. T. Hoyt, L. Vermue, M. Galkin, S. Sharifzadeh, A. Fischer, V. Tresp, and J. Lehmann. Bringing light into the dark: A large-scale evaluation of knowledge graph embedding models under a unified framework. *IEEE Trans. Pattern Anal. Mach. Intell.*, 44(12):8825–8845, 2022.

[5] R. Batra, L. Song, and R. Ramprasad. Emerging materials intelligence ecosystems propelled by machine learning. *Nature Reviews Materials*, 6(8):655–678, Aug. 2021. Number: 8 Publisher: Nature Publishing Group.

[6] D. A. Boiko, R. MacKnight, B. Kline, and G. Gomes. Autonomous chemical research with large language models. *Nature*, 624(7992):570–578, Dec. 2023. Number: 7992 Publisher: Nature Publishing Group.

[7] R. Brighenti, Y. Li, and F. J. Vernerey. Smart Polymers for Advanced Applications: A Mechanical Perspective Review. *Frontiers in Materials*, 7, 2020.

[8] L. Brinson, M. Deagen, W. Chen, J. McCusker, D. Mcguinness, L. Schadler, M. Palmeri, U. F. Ghumman, A. Lin, and B. Hu. Polymer nanocomposite data: Curation, frameworks, access, and potential for discovery and design. *ACS Macro Letters*, 9:1086–1094, 07 2020.

[9] S. Das, S. Sundara, and R. Cyganiak. R2RML: RDB to RDF mapping language. W3C recommendation, W3C, Sept. 2012. https://www.w3.org/TR/2012/REC-r2rml-20120927/.

[10] Data Repository | Materials Genome Initiative. https://www.mgi.gov/infrastructure/data-repository.

[11] Millions of new materials discovered with deep learning. https://deepmind.google/discover/blog/millions-of-new-materials-discovered-with-deep-learning/.

[12] M. Dragoni, G. Rizzo, and M. A. Senese. Chapter 2 - Convology: an ontology for conversational agents in digital health. In S. Jain, V. Jain, and V. E. Balas, editors, *Web Semantics*, pages 7–21. Academic Press, Jan. 2021.

[13] J. H. Dunlap, J. G. Ethier, A. A. Putnam-Neeb, S. Iyer, S.-X. L. Luo, H. Feng, J. A. G. Torres, A. G. Doyle, T. M. Swager, R. A. Vaia, P. Mirau, C. A. Crouse, and L. A. Baldwin. Continuous flow synthesis of pyridinium salts accelerated by multi-objective Bayesian optimization with active learning. *Chemical Science*, 14(30):8061–8069, Aug. 2023. Publisher: The Royal Society of Chemistry.

[14] A massive new effort to name millions sold into bondage during the transatlantic slave trade. https://www.washingtonpost.com/history/2020/12/01/slavery-database-family-genealogy/.

[15] R. Guha and D. Brickley. RDF schema 1.1. W3C recommendation, W3C, Feb. 2014. https://www.w3.org/TR/2014/REC-rdf-schema-20140225/.

[16] M. Harris. 'Room-temperature superconductor' LK-99 fails replication tests. https://physicsworld.com/a/room-temperature-superconductor-lk-99-fails-replication-tests/, Aug. 2023.

[17] K. M. Jablonka, Q. Ai, A. Al-Feghali, S. Badhwar, J. D. Bocarsly, A. M. Bran, S. Bringuier, L. C. Brinson, K. Choudhary, D. Circi, S. Cox, W. A. d. Jong, M. L. Evans, N. Gastellu, J. Genzling, M. V. Gil, A. K. Gupta, Z. Hong, A. Imran, S. Kruschwitz, A. Labarre, J. Lála, T. Liu, S. Ma, S. Majumdar, G. W. Merz, N. Moitessier, E. Moubarak, B. Mouriño, B. Pelkie, M. Pieler, M. C. Ramos, B. Ranković, S. G. Rodriques, J. N. Sanders, P. Schwaller, M. Schwarting, J. Shi, B. Smit, B. E. Smith, J. V. Herck, C. Völker, L. Ward, S. Warren, B. Weiser, S. Zhang, X. Zhang, G. A. Zia, A. Scourtas, K. J. Schmidt, I. Foster, A. D. White, and B. Blaiszik. 14 examples of how LLMs can transform materials science and chemistry: a reflection on a large language model hackathon. *Digital Discovery*, 2(5):1233–1250, Oct. 2023. Publisher: RSC.

[18] K. Janowicz, P. Hitzler, W. Li, D. Rehberger, M. Schildhauer, R. Zhu, C. Shimizu, C. K. Fisher, L. Cai, G. Mai, J. Zalewski, L. Zhou, S. Stephen, S. G. Estrecha, B. D. Mecum, A. Lopez-Carr, A. Schroeder, D. Smith, D. J. Wright, S. Wang, Y. Tian, Z. Liu, M. Shi, A. D'Onofrio, Z. Gu, and K. Currier. Know, know where, knowwheregraph: A densely connected, cross-domain knowledge graph and geo-enrichment service stack for applications in environmental intelligence. *AI Mag.*, 43(1):30–39, 2022.

[19] M. Kerjiwal. *Knowledge Graphs*. MIT Press, 2021.

[20] Materials Project - Home. https://next-gen.materialsproject.org/.

[21] Open-source software enables scientists to expedite research. https://www.afrl.af.mil/News/Article/2775183/open-source-software-enables-scientists-to-expedite-research, Sept. 2021.

[22] OpenRefine. https://openrefine.org/.

[23] OSLNadmin. Polymers take on new importance for Ohio's economy and educators. https://osln.org/2018/11/new-ways-to-learn-about-polymers-launched-as-economic-importance-grows/, Nov. 2018.

[24] M.-E. Papadaki, Y. Tzitzikas, and M. Mountantonakis. A brief survey of methods for analytics over rdf knowledge graphs. *Analytics*, 2(1):55–74, 2023.

[25] Polymer and Responsive Materials Team – Air Force Research Laboratory. https://afresearchlab.com/technology/successstories/polymer-and-responsive-materials-team/.

[26] Polymer Database(PoLyInfo) - DICE :: National Institute for Materials Science. https://polymer.nims.go.jp/.

[27] Polymer Synthesis. https://www.sigmaaldrich.com/US/en/applications/materials-science-and-engineering/polymer-synthesis.

[28] A. Samir, F. H. Ashour, A. A. A. Hakim, and M. Bassyouni. Recent advances in biodegradable polymers for sustainable applications. *npj Materials Degradation*, 6(1):1–28, Aug. 2022. Number: 1 Publisher: Nature Publishing Group.

[29] D. P. Schmidt. *Identifying Knowledge Gaps Using a Graph-based Knowledge Representation*. PhD thesis, Wright State University, 2020.

[30] A. Seaborne and S. Harris. SPARQL 1.1 query language. W3C recommendation, W3C, Mar. 2013. https://www.w3.org/TR/2013/REC-sparql11-query-20130321/.

[31] C. Shimizu, K. Hammar, and P. Hitzler. Modular ontology modeling. *Semantic*, 2021. In Press.

[32] Toolbox Dialogue Initiative – Starting the Dialogue. https://tdi.msu.edu/.

[33] V. Venugopal, S. Pai, and E. Olivetti. MatKG: The Largest Knowledge Graph in Materials Science – Entities, Relations, and Link Prediction through Graph Representation Learning. https://arxiv.org/abs/2210.17340v1, Oct. 2022.