

City Research Online

City, University of London Institutional Repository

Citation: Chen, J., Dong, H., Hastings, J., Jimenez-Ruiz, E., Lopez, V., Monnin, P., Pesquita, C., Škoda, P. & Tamma, V. (2023). Knowledge Graphs for the Life Sciences: Recent Developments, Challenges and Opportunities. Transactions on Graph Data and Knowledge (TGDK),

This is the accepted version of the paper.

This version of the publication may differ from the final published version.

Permanent repository link: https://openaccess.city.ac.uk/id/eprint/31611/

Link to published version:

Copyright: City Research Online aims to make research outputs of City, University of London available to a wider audience. Copyright and Moral Rights remain with the author(s) and/or copyright holders. URLs from City Research Online may be freely distributed and linked to.

Reuse: Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way. City Research Online: <u>http://openaccess.city.ac.uk/</u> <u>publications@city.ac.uk</u>

Knowledge Graphs for the Life Sciences: Recent **Developments, Challenges and Opportunities**

Jiaoyan Chen¹ ⊠ [[]

Department of Computer Science, University of Manchester, Manchester, UK Department of Computer Science, University of Oxford, Oxford, UK

Hang $Dong^2 \square \square$

Department of Computer Science, University of Oxford, Oxford, UK

Janna Hastings³ 🖂 🕒

Institute for Implementation Science in Health Care, University of Zurich, Switzerland School of Medicine, University of St. Gallen, Switzerland

Ernesto Jiménez-Ruiz 🖂 🏠 💿

City, University of London, UK SIRIUS, University of Oslo, Norway

Vanessa Lopez ⊠ IBM Research Europe, Ireland

Pierre Monnin 🖂 🏠 回 Université Côte d'Azur, Inria, CNRS, I3S, France

Catia Pesquita 🖂 LASIGE, Faculdade de Ciências, Universidade de Lisboa, Portugal

Petr Škoda 🖂 回

Department of Software Engineering, Faculty of Mathematics and Physics, Charles University, Prague, Czechia

Valentina Tamma⁴ 🖂 🗅

Department of Computer Science, University of Liverpool, UK

— Abstract -

The term *life sciences* refers to the disciplines that study living organisms and life processes, and include chemistry, biology, medicine, and a range of other related disciplines. Research efforts in life sciences are heavily data-driven, as they produce and consume vast amounts of scientific data, much of which is intrinsically relational and graphstructured.

The volume of data and the complexity of scientific concepts and relations referred to therein promote the application of advanced knowledgedriven technologies for managing and interpreting data, with the ultimate aim to advance scientific discovery.

In this survey and position paper, we discuss

recent developments and advances in the use of graph-based technologies in life sciences and set out a vision for how these technologies will impact these fields into the future. We focus on three broad topics: the construction and management of Knowledge Graphs (KGs), the use of KGs and associated technologies in the discovery of new knowledge, and the use of KGs in artificial intelligence applications to support explanations (explainable AI). We select a few exemplary use cases for each topic, discuss the challenges and open research questions within these topics, and conclude with a perspective and outlook that summarizes the overarching challenges and their potential solutions as a guide for future research.

Authors are listed in alphabetic order with authors' contributions at the end of the article.

Corresponding author

[‡] Corresponding author

[§] Corresponding author

[©] Jiaoyan Chen, Hang Dong, Janna Hastings, Ernesto Jimenez-Ruiz, Vanessa Lopez, Pierre Monnin, Catia etr Škoda, and Valentina Tamma;

⁽cc) licensed under Creative Commons Attribution 4.0 International (CC BY 4.0) Transactions on Graph Data and Knowledge, Vol. 000, Issue 111, Article No. 42, pp. 42:1-42:32

Transactions on Graph Data and Knowledge

GDK Transactions on Graph Data and Knowledge Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

42:2 Knowledge Graphs for Life Sciences

2012 ACM Subject Classification Information systems \rightarrow Graph-based database models, Computing methodologies \rightarrow Knowledge representation and reasoning, Applied computing \rightarrow Life and medical sciences

Keywords and Phrases Knowledge graphs; Life science; Knowledge discovery; Explainable AI Digital Object Identifier 10.1234/0000000.00000000

Funding Jiaoyan Chen: supported by the EPSRC project ConCur (EP/V050869/1).

Hang Dong: supported by the EPSRC project ConCur (EP/V050869/1).

Janna Hastings: supported by the School of Medicine of the University of St. Gallen.

Ernesto Jiménez-Ruiz: supported by the SIRIUS Centre for Scalable Data Access (Research Council of Norway, project 237889).

Catia Pesquita: funded by the FCT through LASIGE Research Unit (ref. UIDB/00408/2020 and ref. UIDP/00408/2020), and also partially supported project 41, HfPT: Health from Portugal, funded by the Portuguese Plano de Recuperação e Resiliência.

Received Date of submission Accepted Date of acceptance Published Date of publishing

Editor Editor Name

1 Introduction

 $_{\rm 2}$ $\,$ The term $life\ sciences\ refers$ to those disciplines that study living organisms and life processes,

and include chemistry, biology, medicine, and a range of other related areas. Research efforts in
 life sciences are increasingly data-driven, as they produce and consume vast amounts of scientific
 data, much of which is intrinsically relational and graph-structured.

⁵ data, much of which is intrinsically relational and graph-structured.

Much of this data is large-scale, complex, and presents many interrelationships and dependencies, thus being well suited to be represented in graph structures. For this reason, graph-based technologies are frequently used in the life sciences, and these disciplines have been drivers and early adopters of innovative methods and associated technologies.

In this brief survey and position paper we discuss recent developments and advances in the use of graph-based technologies in life sciences, and set out a vision for how these technologies will impact these fields in future. We illustrate the contribution in this paper in Figure 1.

We consider Knowledge Graphs (KGs) and their associated technologies to broadly include (*i*) different forms of graph-based representations, (*ii*) the logical languages that assign explicit semantics to such representations, and their associated automated reasoning technologies, and (*iii*) machine learning approaches that ingest data in graph-based representations and that process these graph-based representations to perform some task, *e.g.*, data analytics.

These different forms of graph-based representations can be further categorized based on the type of content represented. We therefore distinguish schema-less and schema-based Knowledge Graphs. More specifically, a typical KG contains either or both a schema part (terminologies or TBox¹) and a data part (facts, assertions, or ABox). The formal semantics of KGs can be expressed with the OWL ontology language².

In the remainder of this paper we will focus on three broad topic areas in which graph-based technologies have been used extensively, and we illustrate each area with some specific projects or use cases that guide our discussion and summary of the challenges that have been encountered.

²⁶ The construction and management of KGs to represent life science knowledge;

²⁷ The use of KGs and associated technologies in the discovery of new knowledge;

¹ We introduce a list of key terms relevant to Knowledge Graphs and Life Sciences in Appendix A.

² Web Ontology Language: https://www.w3.org/OWL/



Figure 1 An overview illustration of definitions (upper right, in grey), topics (left column, in blue), use cases (middle), and challenges (bottom right, in green) for the research of KGs in the life sciences.

The use of KGs in artificial intelligence applications to support explanations (eXplainable AI or XAI).

We then provide a summary of the general challenges across the topics, that include intrinsic 30 characteristics of KGs (e.g., scalability, evolution, heterogeneity) and their operational aspects in 31 the real world (e.g., human interaction, personalization, distributed setting, and representation32 learning). We present the challenges by means of use cases and the current research efforts that 33 address them. It is worth mentioning that while we aim to focus on the life sciences, many of the 34 topics and challenges discussed in this work, especially those of KG construction and management 35 in Section 3, are general and applicable to KGs in other domains such as finance, e-commerce, 36 material, and urban management [114, 32], etc. The KG-based problem modeling and solving 37 approaches in life science knowledge discovery could be applicable for addressing many other use 38 cases and problems in a broader domain of AI for scientific discovery [178, 62]. 39

In the next section, we introduce several different categories of KGs as they have been used in life sciences. Thereafter in Sections 3-5, each of the above topics is described in a dedicated section together with a survey of recent advances. Finally, in Section 6 we synthesize the overarching challenges and trends into a perspective on the outlook for the future.

2 Knowledge Graphs in the Life Sciences

⁴⁵ KGs represent semantically-described real-world entities, typically through ontologies (vocabularies or schemas) [70, 63] and the data instantiating them, and thus provide descriptions of the entities of interest and their interrelations, by means of links to ontology classes describing them, organized in a graph [161]. KGs have been widely adopted in the life sciences, as can be seen in the composition of the Linked Open Data Cloud³, where life sciences represent one of the largest subdomains. A prominent example is the KG representing annotations regarding proteins by means of terms in the Gene Ontology describing different protein functions [4].

⁵² Whilst KGs are becoming increasingly popular in different domains including the life sciences, ⁵³ there is no single accepted definition of KG [45]. A KG can be formally described as a directed,

44

³ http://cas.lod-cloud.net

42:4 Knowledge Graphs for Life Sciences

edge-labeled graph $\mathcal{G} = (V, E)$, where V refers to the vertices or nodes, representing real-world entities of interest (e.g., proteins, genes, compounds, cellular components, but also pathways, biological processes and molecular functions, to name a few) while E refers to the edges in the graph, representing relationships or links between the entities in V (e.g., binds, associates, etc.). These may be represented as statements about entities in the form of RDF⁴ triples: (subject, predicate, object).

However, this formal definition only focuses on the components of KGs, but does not pose any constraint on what a KG should model or represent, and how. This is particularly true in life sciences, where the term *Knowledge Graph* has been used to refer to diverse graph data structures, typically interconnected, but often isolated.

Many of the everyday tasks faced by researchers in this domain require the systematic pro-64 cessing and integration of data and knowledge from data sources that are characterized by het-65 erogeneous syntaxes and structures, formats, entity notation, schemas and scope, e.g., ranging 66 from molecular mechanisms to phenotypes. Researchers in this area have been early adopters 67 of Semantic Web and linked data approaches as a means to facilitate knowledge integration and 68 processing to support tasks including semantic search, clinical decision support, enrichment anal-69 ysis, data annotation and integration. However, a recent analysis of life science open data has 70 identified several stand-alone data sources that exist in isolation, are not interlinked with other 71 sources, and are schema-less (or use unpublished schemas), with limited reuse or mappings to 72 other data sources [90]. Therefore, we can define a life sciences KG, following [133], as a data 73 resource integrating one or more possibly curated sources of information into a graph whose 74 nodes represent entities and edges represent relationships between two entities. This definition is 75 consistent with other definitions found in the literature, e.g., [138]. 76

These considerations underlie the reasons why KGs in life sciences can be of different types, 77 and can be categorized across different dimensions. One of the most critical dimensions (in terms 78 of support for complex queries and integration) is the categorization of KGs into schema-based 79 and schema-less knowledge bases. In turn, the expressivity of the schema provides a further 80 categorization criterion, depending on whether schemas are modelled as simple taxonomies (e.q., 81 the NCBI taxonomy [157] included in the UMLS Metathesaurus [10]), RDFS⁵ vocabularies or 82 (fully axiomatized) OWL ontologies. In particular, this paper refers to this broad definition of 83 KGs, which we then divide into: 84

Schema-less KGs composed of only relational facts in the form of RDF triples. Examples
 include the PharmaGKB dataset, an integrated online knowledge resource capturing how
 genetic variation contributes to variation in drug response [183]. Note that many semantic
 networks (defined in Appendix A) could be assigned to this category as their triples form a
 multi-relational graph.

Schema-based KGs composed of relational facts and their schema (meta information) in *e.g.*,
 RDFS, OWL, and constraint languages such as SHACL⁶. Examples include Wikidata with its
 property constraints, and DBpedia with its DBpedia ontology. Whilst Wikidata and DBpedia
 are general-purpose KGs, they also include large-scale life science knowledge.

Simple ontologies representing taxonomies. Notable examples include the tree structure of
 the UMLS Semantic Network⁷ and the International Classification of Diseases, version 10

96 (ICD-10) [185].

⁴ Resource Description Framework: https://www.w3.org/RDF/

⁵ RDF Schema: https://www.w3.org/TR/rdf-schema/

⁶ https://www.w3.org/TR/shacl/

⁷ https://uts.nlm.nih.gov/uts/umls/semantic-network/root

 Expressive OWL ontologies, with complex axioms beyond simple taxonomies. OWL ontologies 97 may be composed of a TBox and an ABox. Depending on the expressivity of the axioms 98 modeled in the ontology, *i.e.*, the basic statements that an OWL ontology expresses, OWL 99 ontologies can fall into one of the previous categories: for instance, an OWL ontology with 100 just an ABox can be seen as the case above of a KG composed of relational facts alone. In 101 this final category we include fully axiomatized OWL ontologies, e.g., with complex classes 102 and property restrictions. Notable examples of these ontologies include SNOMED CT [39], 103 the Gene Ontology [4, 29], and the Food Ontology (FoodOn)⁸. 104

¹⁰⁵ **3** Knowledge Graph Construction and Management

The adoption of KGs in the life sciences is motivated by the need for standardisation of tax-106 onomies and vocabularies to support the integration, exchange and analysis of data. More re-107 cently, richly annotated data is also being used in combination with machine learning methods 108 for many applications, including helping to overcome issues related to the sparsity of data and 109 helping to select promising candidates for reducing expensive and time-consuming physical ex-110 periments [66]. Graph-based machine learning approaches such as Graph Neural Networks have 111 been applied to a number of life science tasks [51], including drug repurposing [123] and predicting 112 polypharmacy side effects [199]. 113

Given the diverse nature of the knowledge and tasks supported by KGs, the focus of stateof-the-art approaches has been the description of how individual KGs are developed within the specific domain [193], typically in terms of the specific approaches used for the development of the KG (*e.g.*, data extraction process, relation extraction and entity discovery), rather than on the overall development process. More recently, some efforts have focused on providing an overview of development approaches and pipelines for the construction of KGs in the life sciences, and beyond [133, 167]. The process of constructing a KG depends heavily on:

- The type of data sources integrated and annotated by the KG, *e.g.*, CSV files, public and proprietary data sources, structured databases, full-text publications, etc.
- The granularity of the KG to be constructed, e.g., schema-less KG, simple or expressive ontology.
- The usability expectations in downstream applications, e.g., the ability to customize and manipulate the graph to support different use cases, or the ease of consumption as input to machine learning methods [53].

A recent systematic review [167] surveyed different KG development approaches to determine a general development framework. The review identified six main phases that are common across different KG development approaches:

- 131 1) Data source selection.
- 132 2) Ontology construction.
- ¹³³ 3) Knowledge extraction.
- ¹³⁴ 4) Knowledge ingestion and validation.
- ¹³⁵ 5) KG storage and inspection.
- ¹³⁶ 6) KG maintenance and evolution.

⁸ http://foodon.org

42:6 Knowledge Graphs for Life Sciences

¹³⁷ In the remainder of this section we will present the individual phases and the role they play ¹³⁸ in a KG development process by means of two use cases, where we illustrate the construction ¹³⁹ of KGs and discuss how these support knowledge integration and validation (Section 3.2). We ¹⁴⁰ then present some recent technical developments in Section 3.3, while Section 3.4 discusses open ¹⁴¹ challenges for the construction and management of KGs.

¹⁴² 3.1 Knowledge Graph Construction Phases

This section provides more details on the phases involved in the KG construction process, with the aim of identifying recent trends, rather than providing an exhaustive literature survey. These phases are discussed in order of execution, however the *ontology construction* phase can occur either together with the data source selection (if an ontology covering the domain of interest already exists or can be constructed through a set of given requirements) or as part of the *knowledge ingestion and validation* phase, where an ontology is built semi-automatically from the available data or through modularization and alignment of existing ontologies.

150 3.1.1 Data source selection

This phase identifies the data sources that are to be integrated by the KG, which in turn affects 151 the choice of knowledge extraction techniques. Generally, life science KGs ingest knowledge 152 from structured, semi-structured and unstructured data sources. By structured we refer to data 153 modeled according to an existing structure, e.g., data in tables or public or proprietary reference 154 (relational) databases such as UniProt [30] or ChEMBL [52]. Semi-structured data refer to, e.g., 155 XML documents [119], whereas unstructured data refer to data that do not conform to a given 156 structure, *i.e.*, free-text sources, such as scientific publications from PubMed⁹. Data ingested 157 from manually curated databases [133] and semi-structured sources constitute the foundation of 158 a KG [53], generally defining the entities and some of the relations in the KG. This data is then 159 further enriched by performing text mining on large-scale free text sources, in order to extract 160 relationships, which is the objective of the knowledge extraction phase. 161

¹⁶² 3.1.2 Ontology construction

The aim of this phase is to define a common, consensus-based, controlled vocabulary to describe
 the data in an *ontology* [149]. The existence of a common structure, or schema, supports querying,
 integration and reasoning tasks over the KG.

Traditional ontology engineering approaches are divided into top-down or bottom-up. Top-166 down approaches are based on more or less formal ontology engineering methodologies [47, 98, 134] 167 or common practices [3] to build ontologies from a description of the domain elicited from domain 168 experts [132], and/or by reusing or extending existing ontologies [84]. Ontology engineering 169 methodologies define the ontology development process in terms of requirement analysis, entity 170 and property definitions, ontology reuse, validation and population. In contrast, bottom-up 171 approaches utilize semi-automatic data driven techniques, e.g., ontology learning from text [113], 172 and can be used to refine and validate an ontology. These approaches are discussed in more detail 173 when presenting the knowledge ingestion and validation phase. 174

¹⁷⁵ Whilst general purpose ontology engineering methodologies have evolved to be used in the ¹⁷⁶ development of KGs [142], a considerable number of ontologies in the life science domain have been

⁹ https://pubmed.ncbi.nlm.nih.gov

¹⁷⁷ built as part of the Open Biological and Biomedical Ontologies (OBO) Foundry effort,¹⁰ which ¹⁷⁸ defines a set of development principles for biological and biomedical ontologies and provides a suite ¹⁷⁹ of high-quality, interoperable, free and open source tools that support ontology development [118].

3.1.3 Knowledge extraction

Knowledge extraction refers to the identification of entities and their relations from the data 181 sources, which is a crucial step in the development of a KG [167]. Entity extraction identifies 182 entities from the various data sources selected using Natural Language Processing (NLP) ap-183 proaches and text mining techniques to analyse and extract relevant information from large text 184 corpora [181, 106, 73]. Named entity recognition (NER) supports the identification of named 185 entities in text, such as drug names, diseases, or chemical compounds, and their classification 186 according to pre-defined entity types [130]. NER approaches in the life sciences are typically 187 based on labour intensive tasks such as the definition of generic (e.g., orthographic, morpholog-188 ical, or dictionary-based) and specific rules that are typically defined by experts, and are not 189 easily applicable to other corpora [198]. There are a number of issues hindering these approaches: 190 a) the pace of scientific discovery and the identification of new entities; b) the large number of 191 synonyms and term variations associated with an entity; and c) entity identifiers that are com-192 posed of a mixture of letters, symbols and punctuation, often in large sentences [104]. More 193 recent approaches have proposed the use of supervised machine learning methods (e.g., condi-194 tional random fields, or Support Vector Machines, SVMs, neural networks, and neural language 195 models in particular) [115, 88, 36] either in isolation, or combined in hybrid approaches to improve 196 accuracy [152]. 197

Entity recognition generates entities that are isolated and not linked [167]. The goal of Rela-198 tion extraction is to discover relationships of interest between a pair of entities, thus describing 199 their interaction. Relation extraction is a necessary step for entities defined in semi-structured or 200 unstructured sources, whereas structured data sources are characterized by explicitly identifiable 201 relationships. Typical approaches for relation extraction include rule-based [77, 148, 147], super-202 vised [109, 50] and unsupervised approaches [101, 133]. Rule-based relation extraction identifies 203 keywords (based on existing ontologies or expert defined dictionaries) and grammatical patterns 204 to discover relations between entities. Supervised relationship extraction methods utilize publicly 205 available pre-labelled datasets (e.g., BioInfer [144] or BioCreative II [100]) to construct generalized 206 patterns that separate positive examples (sentences implying the existence of a relationship) from 207 negative ones. Supervised approaches include SVMs, Recurrent Neural Networks (RNNs) and 208 Convolutional Neural Networks (CNNs) [7, 133]. Unsupervised relation extraction methods [116] 209 have emerged to address the lack of scalability of supervised relation extraction methods, due to 210 the high cost of human annotation. Unsupervised methods involve some form of clustering or 211 statistical computation to detect the co-occurrence of two entities in the same text [133]. 212

More recently, end-to-end approaches (End-to-End Relation Extraction – RE) have been used to tackle both tasks simultaneously. In this scenario, a model is trained simultaneously on both the NER and Relation Extraction objectives [76]. Furthermore, rule-based approaches can be combined with relation classification using specialized pre-trained language models adapted for life science domains, *e.g.*, BioBERT [105], SapBERT [111], and RoBERTa-PM [107], to name a few. There is also a recent trend to probe and prompt pre-trained language models to extract relations (*e.g.*, disease-to-disease, disease-to-symptoms) [190, 166].

¹⁰ https://obofoundry.org

220 3.1.4 Knowledge ingestion and validation

The aim of this phase is to ingest the entities and relationships extracted in a previous phase, which models knowledge from different sources. These entities and relations can be incomplete, ambiguous or redundant, and need to be appropriately aligned and integrated, and finally annotated according to the ontology constructed in phase 2.

Knowledge integration or fusion can critically improve the quality of data by performing *entity* 225 resolution, i.e., the detection of different descriptions of the same real-world entity (also called 226 entity matching, deduplication, entity linkage or entity canonicalization), prior to ingesting them 227 in the KG. This reconciliation step is particularly crucial in the life sciences, where duplication can 228 be caused by data modelled using different vocabularies or ontologies, or when data is extracted 229 from literature sources that are rapidly changing. The severity of the ambiguity depends on the 230 number of ontologies available for the domain. For instance, the number of gene vocabularies 231 is far smaller than the number of disease vocabularies that could be present in the ingested 232 datasets. Linking these entities requires costly alignment processing; in particular the alignment 233 of disease entities is especially problematic given the number of different coding systems, whose 234 conversion is often not trivial [53]. We further explore this issue in two of the use cases presented 235 in Section 3.3, where we explore the problem of aligning vocabularies and ontologies through the 236 use of mapping repositories and instance matching in automated clinical coding. 237

Entities are assigned unique identifiers (URI or IRI) that support the definition of bespoke namespaces, and support integration by reusing identifiers in related namespaces. Entity resolution is based on clustering similar entities together in a *block*, where similarity measures are used to detect duplicates [167]. Typical methods include sorted neighborhoods and traditional blocking; and machine learning methods are commonly used for similarity computation, *e.g.*, feature vector computation [96].

This phase may also include the bottom-up construction of the ontology for those applications where a top-down approach is not feasible. Bottom-up approaches extract the relevant knowledge first, and then they construct the data schema / ontology based on the extracted data, typically using (semi-)automated methods, based on machine learning. Ontologies define the structure of the knowledge graph, which supports querying and data analytics. In bottom-up ontology development the structure of the knowledge graph is determined based on the extracted knowledge, thus providing a structure for this knowledge [71].

Often the construction of ontologies (either bottom-up or top-down) relies on the ability to correctly align and reuse entities defined across different domains and KGs. Furthermore, reuse of (or conformance to) existing upper level ontologies, *e.g.*, BFO (Basic Formal Ontology) [3] provides the basis for the consistent and unambiguous formal definition of entities and relations that prevents errors in coding and annotation. The alignment of ontologies in life sciences and other domains is an active area of research, and we provide an overview of recent technical developments and challenges in Section 3.3.

Whilst bottom-up approaches, especially those based on alignment, are becoming more viable, especially given the support of language models, such as BERT [65], their performance is not always adequate for the task, as discussed in the second challenge in Section 3.4.

Knowledge enrichment and completion improve the KG quality by performing reasoning (KG materialization), inference [58] and optimization. Reasoning and inference support the assertion of new relations based either on logical reasoning (*e.g.*, [131, 173]) or machine learning techniques (*e.g.*, statistical relational learning or through embedding based link predictors for new concepts [35, 36, 68, 78] and node classifiers, also called KG refinement [138]). The extent and type of logical inferences depends on the expressivity of the ontology built in phase 2, or in a bottom-up fashion in this phase, together with any associated mappings. Description Logic for-

42:9

malisms, such as OWL, use logic-based reasoning for detecting and correcting incorrect assertions
 and ontology alignments [25].

270 3.1.5 KG storage and inspection

KGs need to be accessible to support a variety of different tasks, beyond the mere integration of 271 different knowledge sources, and thus KG storage management [167, 145, 180] is an active area 272 of research. Current KG storage mechanisms are divided into relation based stores (e.q., [1]) and 273 native graph stores (e.g., [200]). Relational KG stores, either based on relational databases or 274 through NOSQL databases and / or triple stores such as Jena TDB¹¹, have reached a considerable 275 level of maturity and have been optimized in order to avoid common problems, e.g., a large number 276 of null values in columns or optimized query performance [145]. Graph databases store nodes, 277 edges and properties of graphs natively, and support query and graph mining tasks. Examples 278 of state of the art implementations include Neo $4J^{12}$, GraphDB¹³, and RDFox¹⁴. The evolution 279 of the performance of these systems has been the object of systematic studies [9], whereas [171] 280 explicitly focuses on biomedical use cases. 281

Storage management has implications on the ways KGs support expressive queries for nodes and edges and visualization, to support data analysis, navigation and discovery of related knowledge [96, 165]. Graph databases often provide built-in tools for visualization, *e.g.*, Neo4J, whereas different Javascript libraries (*e.g.*, SigmaJS¹⁵) are available for developing visualization front ends. Support for complex queries is also either built in a graph database or a triple store by supporting the SPARQL query language [143, 200], or proprietary query languages such as Cypher [49], supported by Neo4J.

²⁸⁹ 3.1.6 Knowledge maintenance and evolution

Given the rapid scientific development in the life sciences, and the consequent continuous update of ontologies for this domain, artefacts annotated with these ontologies can become outdated very quickly, and require some form of update (also called ontology extension). These update mechanisms need to be automated to ensure that they scale to the size of KGs. Automatic update approaches are based on the periodical detection and extraction of new knowledge that is then mapped to existing entities and relations in the KG [186].

Update mechanisms are typically based on the detection of *changes* [124] that can affect an 296 ontology, e.g., addition, removal or modification of meta-entities (*i.e.*, entities, relations and their 297 definitions). These changes include renaming concepts and properties, setting domain and range 298 restrictions, or setting a subsumption relation. To date, the most comprehensive account of 200 ontology change is given in [48], where change is described for different sub-fields, e.g., ontology 300 alignment, matching and mapping, morphisms, articulation, translation, evolution, debugging, 301 versioning, integration and merging; each with different requirements and implications. The 302 study [140] further investigates the impact of biomedical ontology evolution on materialization. 303

Currently available tools and methodologies use (semi)-automated methods to perform many of the operations that trigger a change in an ontology and the consequent creation of a new version [56, 65]. Different ontology management platforms and portals mandate different principles

 $^{^{11} \}tt https://jena.apache.org/documentation/tdb/index.html$

¹² https://neo4j.com

¹³https://graphdb.ontotext.com

¹⁴ https://www.oxfordsemantic.tech/product

¹⁵ https://github.com/jacomyal/sigma.js

42:10 Knowledge Graphs for Life Sciences

and frameworks for handling ontology versioning (*e.g.*, OBO foundry¹⁶ or BioPortal¹⁷), but these are typically implemented by ontology developers with limited tool support. Section 3.3 presents an example of automated ontology extension that relies on machine learning to cope with the scale of data.

311 3.2 Examples of Life Science KG Construction

In this section we provide two examples of life science KGs that illustrate in practice the phases composing the generic KG construction process discussed in Section 3; namely a KG for Pharmacogenomics, PGxLOD [121], and one for Ecotoxicological Analysis, TERA [127, 128].

Alignment for Knowledge Validation: An Example of Pharmacogenomics. As men-315 tioned in Section 3, the task of aligning knowledge in KGs supports several downstream appli-316 cations and domains. For instance, pharmacogenomics studies the influence of genetic factors 317 on drug response phenotypes (e.g., expected effect, side effect). Hence, pharmacogenomics is of 318 interest for personalized medicine. The atomic knowledge unit in pharmacogenomics is a ternary 319 relationship between a drug, a genetic factor, and a phenotype. Such a relationship states that 320 a patient being treated with the specified drug while having the specified genetic factor may 321 experience the described phenotype. Semantic Web and KG technologies have been employed in 322 this application domain, for example by building ontologies in which patients and pharmacoge-323 nomic knowledge are represented, and then using deductive reasoning mechanism to conditionally 324 recommend genetic testing before drug prescription [156]. However, the knowledge relevant to 325 pharmacogenomics is scattered across several sources including reference databases such as Phar-326 mGKB, and the biomedical literature. Additionally, this knowledge may lack sufficient validation 327 to be implemented in clinical practice. For example, some relationships may have only been 328 observed in smaller cohorts of patients or in non-replicated studies. Hence, there is a need to 320 align different sources of pharmacogenomic knowledge to detect additional evidence validating 330 (or moderating) a knowledge unit. To this aim, the PGxLOD KG was proposed [121]. Automatic 331 knowledge extraction approaches were applied on semi-structured and unstructured data from 332 PharmGKB and the biomedical literature to represent their knowledge in the KG. Then, match-333 ing approaches were developed to align knowledge units from various sources [120, 122]. The 334 resulting alignments outlined some agreements between PharmGKB and the biomedical litera-335 ture, which was expected since PharmGKB is manually completed by experts after reviewing the 336 literature. Interestingly, this automatic knowledge extraction pipeline could guide the manual re-337 view process by automatically pointing out studies confirming or mentioning a pharmacogenomic 338 knowledge unit. 339

Knowledge Integration: An Example of Ecotoxicological Analysis. In ecotoxicological 340 analysis, data and knowledge from different domains such as chemistry and biology are often 341 needed. These are usually located in different sources such as spreadsheets or CSV files for 342 local experimental results, open databases for public research results, and ontologies for domain 343 knowledge. Thus knowledge integration becomes a critical and fundamental challenge before 344 real analysis can be conducted. In the study by Myklebust et al. [127, 128], which aims to 345 predict adverse biological effects of chemicals on species, a toxicological effect and risk assessment 346 KG named TERA was constructed for knowledge integration. TERA includes three sub-KGs: 347 (i) the Chemical sub-KG, which is constructed by integrating the vocabulary MeSH (Medical 348

 $^{^{16}}$ http://www.obofoundry.org/principles/fp-004-versioning.html

¹⁷ https://bioportal.bioontology.org

Subject Headings) with selective knowledge from two chemical databases PubChem and ChEMBL 349 utilizing the chemical mappings in Wikidata; (ii) the Taxonomy sub-KG, which is constructed by 350 integrating EOL (Environment Ontology for Livestock) and the NCBITaxon ontology utilizing 351 NIBI-EOL mappings in Wikidata; and (iii) the ECOTOX sub-KG, which is composed of RDF 352 triples transformed from experimental risk results and is aligned with the other two sub-KGs by 353 the ontology alignment system LogMap [82] and the chemical mappings in Wikidata. Another 354 example of knowledge integration is for drug repurposing, where the KG Hetionet¹⁸ is created by 355 integrating 29 public resources, including biomedical KGs and other types of data [69]. 356

357 3.3 What has been done: recent technical developments

Given the many existing ontologies in life sciences, *e.g.*, ontologies available in the OBO Foundry collection or in BioPortal [135], KG construction usually involves the reuse, alignment, and enrichment of state-of-the-art ontologies. The existing ontologies in life sciences need to be updated given the new discoveries in the field. This is broadly a key issue in the management, maintenance, and evolution of ontologies. We select a few promising use cases below to highlight some recent developments that support the KG construction in the life sciences.

Repositories of Ontologies and Mappings. Ontologies and their mappings play a central role 364 in semantically enabled products and services consumed by life science companies, academic in-365 stitutions and universities, as highlighted by the Pistoia Alliance ontology mapping project [60].¹⁹ 366 Ontology mappings are essential in knowledge graph construction tasks to bridge the knowledge 367 provided by different ontologies and expand their coverage. Ontology mappings can also play a 368 key role when identifying the right ontologies to be reused as they will enable the retrieval of 369 the relevant (overlapping) ontologies for the domain of interest. For this reason, a number of 370 notable efforts in life sciences have created large repositories of ontologies and mappings to serve 371 the research within the community. Prominent examples include the UMLS Metathesaurus [10], 372 BioPortal [135, 155], MONDO [175], and the EBI services: OLS [177], OXO [86] and the RDF 373 platform [87]. The UMLS Metathesaurus is a comprehensive effort for integrating biomedical 374 ontologies through mappings. In its 2023AA version, it integrates more than two hundred vo-375 cabularies, with more than 3 million unique concepts and more than 15 million concept names. 376 BioPortal is a repository containing more than 1,000 biomedical ontologies and more than 79 377 million lexically computed mappings among them (as of July 13, 2023). The Mondo Disease 378 Ontology (MONDO) is a manually curated effort to harmonize and integrate disease concep-379 tualizations and definitions across state-of-the-art ontologies (e.g., HPO [99], DO [158], ICD, 380 SNOMED CT, etc.). The services provided by the European Bioinformatics Institute (EBI) also 381 deserve a special mention. The Ontology Lookup Service (OLS) has become a reference to explore 382 the latest versions of more than two hundred ontologies via its graphical interface or program-383 matically via its API. OxO is a repository of ontology mappings and cross-references extracted 384 from the OLS and UMLS. OxO allows users to visually traverse the graph of mappings to identify 385 additional potential mappings beyond direct ones (*i.e.*, multi-hop mappings). Finally, the EBI 386 RDF platform provides a unified KG with all the RDF resources at the EBI. Complementary to 38 the efforts from the life sciences, the Semantic Web has also contributed to the systematic eval-388 uation of mappings in public repositories (e.g., [83, 46]) and mappings produced by automated 389 ontology mapping systems (e.g., the Ontology Alignment Evaluation Initiative (OAEI) [141]). 390 Automatically generated mappings of high quality have the potential to be integrated within the 391

¹⁸ https://github.com/hetio/hetionet

¹⁹ https://www.pistoiaalliance.org/projects/current-projects/ontologies-mapping/

42:12 Knowledge Graphs for Life Sciences

³⁹² aforementioned repositories and hence, the OAEI has always had a special focus on life science ³⁹³ test cases with evaluation tracks like Anatomy [41], LargeBio [85], Phenotype [61] and the newly ³⁹⁴ created track BioML [66]. The Simple Standard for Sharing Ontological Mappings (SSSOM) [117] ³⁹⁵ represents a joint effort between the life sciences and Semantic Web communities to facilitate the ³⁹⁶ exchange of mappings across different parties and repositories, while keeping the provenance and

³⁹⁷ other relevant characteristics of the mappings.

Ontology Extension. Ontology extension in life sciences aims to connect new concepts and 398 their relations to an ontology from updated sources, e.g., scientific papers in PubMed and chemical 399 information in PubChem²⁰. Manual ontology extension, while essential for the development of 400 gold standard resources, is not scalable to the full scope of large domains due to its high cost and 401 low efficiency, and sometimes is even unfeasible as human beings may not be able to review the 402 quantities of new information at the rate they become available. Thus machine-learning-based, 403 automated methods are needed. One recent example is the use of deep learning, specifically 404 a Transformer-based model, to categorize new chemical entities within the ChEBI ontology²¹ 405 [55]. In addition, recent studies have explored enriching SNOMED CT by mining new concepts 406 from texts [36] and placing them into the ontology [112, 35]. A new concept can be identified by NIL entity linking, *i.e.*, exploring unlinkable mentions, usually through setting a "linkable" 408 score threshold or through classification [36]. Resolution and disambiguation of NIL mentions 409 with clustering can help to represent NIL entities [68, 94]. For concept placement, similar to 410 the aforementioned CHEBI ontology extension [55], machine learning, especially in the form of 411 Transformer-based deep learning, has been applied to predict subsumption relations between 412 a new concept and the existing concepts. Complex concepts in OWL ontologies that contain 413 logical operators (e.g., existential quantifier and conjunction in SNOMED CT) can be supported 414 in subsumption prediction [24] and new concept placement [35]. Another group of studies use 415 post-coordination or formalising a new term with existing concepts and attributes [17, 95], which 416 is similar to composing subsumption axioms with complex concepts. The methods include using 417 lexical features [95], word embeddings and KG embeddings [17]. Pre-trained and Large Language 418 Models, through fine-tuning, zero-shot and few-shot prompting have the potential to support the 419 mining [36] and placement of new concepts (e.g., by subsumption prediction [24, 67]). 420

Instance Matching: Automated Clinical Coding. A main source for patients' KG construc-421 tion is Electrical Health Records (EHR). Using medical ontologies as backbones, it is possible to 422 add a layer of data by instance matching (or patient matching) through *Clinical Coding*. Clinical 423 coding is the task of transforming medical information in EHR into structured codes described 424 in medical ontologies [37], e.g., ICD and SNOMED CT. Recent approaches mainly formulate the 425 problem as a multi-label classification problem. Various neural network architectures have been 426 proposed and knowledge plays a key role to enhance the neural architectures [37, 81]. Pre-trained 427 language models, e.g., BERT [33], have been applied to clinical coding and gradually achieved 428 better results with adapted modelling methods and more advanced language models, e.g., PLM-429 ICD [72] with RoBERTa-PM [107], according to studies [37, 44, 80]. Other studies formulate the 430 task as a Named Entity Recognition and Linking (NER+L) problem, by extraction of concepts 431 and linking them with the ontologies [37]. Overall, the recent progress in clinical coding, along 432 with the advent of Large Language Models (LLMs) suggests a trend in this area for patients' KG 433 construction from EHR. However, there is still room for improvement in knowledge integration to 434 better address explainability (see Section 5 for more details) and in zero-shot learning problems, 435

 $^{^{20}} https://pubchem.ncbi.nlm.nih.gov/$

²¹ https://www.ebi.ac.uk/chebi/

i.e., for classifying into rare codes or concepts [37, 44, 81]. There are also further recent examples of instance matching with EHR data, including the works [16, 169].

3.4 What are the challenges?

KG construction and management often play a fundamental role in supporting life sciences with
computation. There are still quite a few technical challenges, and many of the current tools and
algorithms can be improved by modern machine learning and AI techniques. Here we present
some critical and fundamental technical challenges.

- How to construct a customized KG? For a specific application, we often need to extract 443 relevant data and knowledge from multiple sources, and at the same time integrate extracted 444 knowledge from different sources. Considering a case study of personal health assistance, 445 a customized KG with knowledge of at least exercise (sports), food, disease and medicine 446 are required, while fine-grained knowledge of these aspects will lie in different domain KGs. 447 The key challenge for integrating different ontology modules lies in estimating the seman-448 tic similarity and discovering the equivalence of two knowledge elements with their contexts 449 considered, as well as the subsequent refinement like KG completion and knowledge represen-450 tation canonicalization. Adequate tool support to minimize manual curation but enabling the 451 user involvement when required is also paramount (e.g., [108]). 452
- How to ensure adequate performance using machine learning based approaches for 453 automated KG construction? At the TBox level, the state-of-the-art alignment between 454 classes (especially for subsumption relations) seems to not yet be achieving good enough 455 performance, as reflected in recent biomedical ontology alignment benchmarking [66]. At the 456 ABox level, predicting missing facts for practical KG construction expects high precision (e.g., 457 beyond 90% or 95%) but only a few relations can be populated with a precision above 80%458 using prompt learning with BERT as evaluated in [176]. This is also the case to associate 459 patients' EHR (as a part of ABox) with clinical codes or concepts in medical ontologies, 460 where a micro F_1 score is below 60% [37]. Learning subsymbolic representations (see defined 461 in Appendix A) of KG and data sources may help address the challenge. Transformer-based 462 language models have achieved great performance in recent years. Among them, pre-trained 463 language models such as BERT have been applied for KG construction with a promising 464 performance achieved (see e.g., the package DeepOnto [65]), while the more recent and more 465 powerful generative language models like GPT series [14] have not been well applied at the 466 time of writing, especially in the life science domain. 467
- How to ensure reliable semi-automated deep learning-based KG construction with 468 human interaction? Many tasks in the KG life cycle unavoidably rely on human experts 469 to achieve consensus on reliable knowledge; on the other hand, as the automated KG con-470 struction process is growing opaque with deep learning methods, it is important to ensure 471 trustworthiness and reliability [194]. Apart from enhancing performance metrics with novel 472 methods, results with certain explainability are needed, for example, highlighting key parts in 473 the data input when they are used as sources for KG construction. We discuss other aspects of 474 explainability with KG, on life science knowledge discovery and healthcare decision making, in 475 Section 5. Human-in-the-loop learning design for explainable KG construction may ensure the 476 use of experts' knowledge for the task across the KG life cycle, which still remains a challenge 477 for future research [194]. 478

479 **4** Life Science Knowledge Discovery

Research into AI technologies – including machine learning and KG-based reasoning – to accelerate the pace of scientific discovery is an emerging and rapidly developing field. The challenge lies
in assisting scientists to uncover new knowledge and solutions, such as discovering novel therapeutic opportunities, identifying candidate molecular drugs to treat complex diseases or alternatively
new uses for existing drugs, and supporting more personalized predictions.

Knowledge Graphs are powerful tools for representing complex biomedical knowledge, including molecular interactions, signalling pathways, disease co-morbidities, and more. Overviews 486 of graph representation learning in biomedicine for healthcare applications and polypharmacy 487 tasks are presented in [110] and [54] respectively. In graph representation learning, the graph's topology is leveraged to create compact vector embeddings. Through nonlinear transformations, 489 high-dimensional information about a node's graph neighborhood is distilled into low-dimensional 490 vectors, where similar nodes are embedded close together in the vectorial space. Embeddings have 49: been shown to be valuable for handling numerous relations in a KG while efficiently exploiting re-492 lation sparsity using vector computations. These optimized representations are subsequently used 493 to train downstream models for various tasks, such as predicting property values of specific nodes 494 (e.q., protein function), predicting links between nodes (e.q., binding affinity between molecules)495 and protein targets), or performing classification tasks (e.g., predicting the toxicity profile of a 496 candidate drug, or risk of readmission for a patient). 497

It is worth mentioning that among the existing works for life science knowledge discovery, 498 different kinds of KGs have been exploited. The schema-less KG can be used to model different 499 kinds of interaction between instances such as proteins and drugs; the taxonomy alike simple 500 ontology is often used to represent concepts and their hierarchy such as protein functions defined 501 in the gene ontology, chemical compounds, species, and diseases; expressive OWL ontologies 502 and schema-based KGs can be used to model complex logical relationships between concepts, 503 besides simple interaction between instances. Such diverse knowledge representation capabilities 504 make KGs more flexible in modeling the input data and prediction targets of different knowledge 505 discovery tasks, than graphs and tabular data that are widely used in previous pure machine 506 learning-based methods. 507

In the following, we present some typical use cases, where machine learning techniques (including graph representation learning and language models) are applied over KGs built from diverse sources and domain ontologies, to facilitate life science discovery.

4.1 What has been done: use cases and their recent developments

Therapeutics and Drug Discovery: Learning a representation using multi-modal and 512 heterogeneous knowledge. Drug discovery entails exploring an extremely large space of poten-513 tial drug candidates. AI can help to accelerate this process by narrowing down the most promis-514 ing candidates before expensive experimentation. The key to leveraging predictive and generative 515 models for candidate solution generation lies in learning an effective multi-modal representation 516 of protein targets, molecules and diseases among others. Recent research has focused on applying 517 language models over large databases of proteins or molecules for self-supervised representation 518 learning, such as ESM [151] and ProteinBERT [11] for protein sequences, or Molformer for the 519 molecule simplified molecular-input line-entry system (SMILES) [154]. These models have exhib-520 ited remarkable success in tasks such as predicting protein interactions, binding affinity between 521 drugs and targets, and protein functions and structures. However, these existing pre-trained 522 sequence-based models often neglect to incorporate background knowledge from diverse sources, 523 for example, biological structural knowledge. 524

Nonetheless, recent research indicates that incorporating existing expressive factual knowl-525 edge can improve results in downstream machine learning tasks. To enhance Protein Language 526 Models (PLM), approaches such as OntoProtein [195] and KeAP [197] use a KG of protein se-527 quences augmented with textual annotations from the Gene Ontology (GO). OntoProtein was 528 the first to inject gene ontology descriptions into a PLM for sequences to predict protein interac-529 tions, function and contact prediction. OntoProtein proposes to reconstruct masked amino acids 530 while minimizing the embedding distance between the contextual representation of proteins and 531 associated knowledge terms. Similarly, ProtST [189] uses a dataset of protein sequences aug-532 mented with textual property descriptions from biomedical texts and jointly trains a PLM with 533 a biomedical language model. 534

Knowledge Graphs are suitable data models for expressing heterogeneous knowledge and fa-535 cilitating end-to-end learning [184]. An entity in a KG can have multiple attributes with different 536 modalities - where each modality provides extra information about the entity - as well as relations 537 to and from entities in other sources. Graph Neural Networks (GNN) have been used to capture 538 inter-dependencies and diverse types of interactions between heterogeneous entity types and mul-539 timodal attributes in KGs [103]. They achieve this by iteratively aggregating information from 540 neighbouring nodes (through a process called message passing) and employing scoring functions 541 to optimize the learned embeddings for downstream tasks. Otter-Knowledge [103] incorporates a 542 heterogeneous KG (schema-based, containing concepts and their attributes) from diverse sources 543 and modalities, *i.e.*, each node has a particular mode that qualifies its type (text, image, protein 544 sequence, molecule, etc.) and initial embeddings for each node are computed based on their 545 modality. A GNN is then used to enrich protein and molecule representations and train a model 546 to produce final node embeddings. The model is able to produce representations for entities 547 that were not seen during training and achieve state-of-the-art results in the Therapeutic Data 548 Commons (TDC) benchmarks [75] for drug-target binding affinity prediction. TxGNN [74] uses a 549 GNN pre-trained on a large heterogeneous, multi-relational KG of diseases and therapeutic can-550 didates constructed from various knowledge bases. TxGNN obtains a signature vector for each 551 disease based on its neighboring proteins, exposure and other biomedical entities to compute a 552 disease similarity and predict drug indication/contraindication for poorly characterized diseases. 553

Protein Function Prediction with the Gene Ontology. Conducting physical experiments 554 for identifying protein functions is time and resource consuming. With the development of ma-555 chine learning, protein function prediction (which is the task of predicting a given protein with 556 multiple and potentially hierarchical classes – functions – defined in GO) has been widely inves-557 tigated in recent years [196, 174]. A large part of these works such as GOLabler [192] focus on 558 exploring feature extraction, feature ensemble, and automatic feature learning of the proteins. 559 For example, GOLabler [192] utilizes five kinds of different protein sequence information while 560 DeepGraphGO [191] builds a network of proteins and learns protein features via a Graph Neural 561 Network. Recent methods attempt to further exploit inter-function (class) relationships that are 562 defined in GO for better performance. For example, DeepGOZero [102] and HMI [188] use formal 563 semantics including the class hierarchy, class disjointness axioms and complex class restrictions 56 in OWL as additional constraints for training the multi-label classifier for protein function pre-565 diction. Protein function prediction is a representative multi-label classification problem where 566 complex relationships of the labels are defined in a KG and can be used for performance aug-567 mentation. It is quite common in machine learning applications in the life sciences, such as the 568 569 above mentioned automated clinical coding where the codes' semantics are modeled by the ICD ontology, and ecotoxicological effect prediction where the multiple affected species of a chemical 570 to predict form a taxonomy. 571

572 Predictions for Healthcare using Ontologies with Clinical Data. Digital Healthcare

42:16 Knowledge Graphs for Life Sciences

involves predictions using clinical data and ontologies, including diagnosis (e.g., rare diseases) 573 and procedure predictions (e.g., ICU readmissions). A related concept is personalized medicine, 574 which is achieved through the matching and fusion of knowledge from diverse sources, and plays 575 a significant role in the prediction tasks. This often involves matching multiple ontologies [159], 576 integrating curated databases (e.g., pharmacogenomics, molecules and proteins knowledge bases), 577 mining knowledge from scientific literature [187] and person-centered clinical knowledge extracted 578 from EHR or claim data, with distinguishing risk factors or cohorts' demographics (e.g., age and 579 gender), which could enhance predictions related to adverse effects [126] or rare diseases for 580 which there are not enough labeled datasets [2]. For example, SHEPHERD [2] incorporates a 58 multi-relational KG (extracted from PrimeKG [20]) of diseases, phenotypes and genes, and lever-582 ages patient simulated data to discover novel connections between patients' clinical, phenotype 583 and gene information to accelerate the diagnoses of rare diseases. Knowledge-guided learning 584 is achieved by training a GNN to represent each patient's subgraphs of phenotypes in relation 585 to other gene, phenotype, and disease associations within the KG, such that embeddings are 586 informed by all of the existing biomedical knowledge captured in the network topology. 587

The approach in [16] constructs a KG (using expressive OWL ontologies) to predict ICU 588 (intensive care units) readmission risk by enriching EHR data with semantic annotations from 589 various biomedical ontologies in BioPortal. These predictions are based on KG embedding, such 590 as RDF2vec, OPA2vec, and TransE, and classical machine learning methods, such as Logistic 591 Regression, Random Forest, Naive Bayes and Support Vector Machines. Drawing from the Health 592 & Social Person-centric Ontology (HSPO) [168], which focuses on multiple clinical, social and 593 demographic facets for a patient or cohort, the approach presented in [169] builds a person-594 centric KG (expressive OWL ontology with TBox and ABox) from structured and unstructured 595 data in EHR). Subsequently, a representation learning approach using GNNs is used to predict 596 readmissions to the ICU. 597

⁵⁹⁸ 4.2 What are the challenges?

⁵⁹⁹ We present four of the open challenges to unlock the full potential of methods to advance knowl-⁶⁰⁰ edge discovery for the life sciences using KGs, based on the use cases above.

How to incorporate the semantics from a KG in machine learning? Many life 601 science knowledge discovery tasks are modeled as a machine learning classification problem, 602 whose input and output labels have additional valuable information in one or multiple ex-603 ternal KGs. The challenge lies in extracting this information, optionally encoding it into 60 vector representations, and injecting that knowledge into machine learning and pre-trained 605 language models. Doing this effectively remains an important open challenge especially for 606 protein-related pre-trained language models [195, 189, 197]. Besides improving the accuracy 607 in knowledge discovery, injecting semantics from KGs can also contribute to making the model 608 more explainable (see Section 5), but to this end, much research is still required. 609

How to deal with the long-tail phenomenon in machine learning with KGs? In 610 machine learning classification for real-world life science knowledge discovery, the candidate 611 labels often exhibit a long-tailed distribution, *i.e.*, a small ratio of them are common with a 612 large number of training samples available, while most of them are infrequent or even have 613 never appeared before. For example, imbalance in training data may occur for rare diseases 614 or adverse drug effects that affect only a small portion of the population [2, 74, 38]. KGs 615 sometimes have encoded the relationships of the labels, and could be used to help train the 616 model for predicting those long-tailed labels or enable the inference of such labels. 617

How to create an efficient multi-modal representation of knowledge to enable dis-618 covery? Most current state-of-the-art methods build learned graph representations based on 619 isolated modalities. Multimodal KGs can explicitly capture labelled nodes and edges, each 620 with well-defined meanings, across heterogeneous node types, relations and modalities (such 621 as text, images, protein sequences, molecules fingerprints, diseases and more) [20, 103]. Incor-622 porating KGs with multiple modalities for representation learning requires computationally 623 scalable methods to compute the initial embeddings for each modality, as a preliminary step 624 to learn computable representations of large knowledge. Furthermore, robust learning tech-625 niques are needed for generalizing the learned representations to nodes with unseen or missing 626 modalities, thereby enabling the discovery of new knowledge. An example would be inferring 627 properties of proteins for which only the sequence is known. 628

How to efficiently utilize and fuse heterogeneous datasets, such as human-curated 629 domain knowledge bases, scientific literature and person-centered health records, 630 for knowledge discovery? State of the art shows that representations can be enhanced by 631 incorporating richer information available across different sources [74, 103, 159]. Bringing in 632 more data during training is needed to learn representations that can be applied to a broader 633 range of downstream prediction tasks. However, learning from large and diverse KGs requires 634 addressing challenges such as alignment, noise handling, balancing rich expressive knowledge 635 with scalability and dealing with knowledge inconsistency. Moreover, more robust learning 636 methods are needed for generalizing the learned representation to multiple downstream tasks 637 (e.g., knowledge-aware transfer, zero-shot and few-shot learning [23]). An important aspect 638 in this regard is addressing the disparity between all of the knowledge accessible during pre-639 training and the knowledge accessible or relevant for the downstream fine-tuning [74, 103]. 640

5 Knowledge Graphs for Explainable AI

Machine Learning (ML) and Artificial Intelligence (AI) methods are widely employed to tackle 642 complex problems in many domains, including life sciences such as chemistry or biomedicine. Yet 643 many of those methods operate as a "black-box", not enabling domain experts to understand 644 the reasoning behind their predictions [93]. This is a major concern, especially for applications 645 in areas with a potential impact on human lives, or areas with legally enforced accountability 646 or transparency [146]. Moreover, understanding the workings of AI methods is also crucial in 647 the context of scientific applications, such as those described in Section 4, where explaining the 648 prediction process can help elucidate natural phenomena [42]. 649

One way to address this issue is to employ the methods of eXplainable Artificial Intelligence 650 (XAI). Although this is a topic long explored in the AI research community, there is still no 651 widely-accepted definition of explainability, with many terms being used interchangeably, such as 652 interpretability, comprehensibility, understandability and transparency [8]. Barredo et al. define 653 explainability as the ability of a model to make its functioning clearer to an audience [8]. A 654 slightly different definition is given in the previous survey [57]: "an interface between humans 655 and a decision maker that is at the same time both an accurate proxy of the decision maker 656 and comprehensible to humans". Both definitions focus on the audience, for whom is the model 657 explainable, but the second suggests an explanation is another artefact produced by a model or 658 alongside the model. 659

There are two distinguishable audiences in the context of the life sciences: scientists (researchers) and healthcare practitioners [170]. For the first group, the explanation is used as a guide to understanding within life sciences research for scientific discovery. As a result, the explanation may exist in a well-bounded context of a hypothesis or research project. On the other

42:18 Knowledge Graphs for Life Sciences

hand, practitioners are involved directly in decisions with impact on healthcare. They need to
 consider the output of the model in an open context, and sometimes also to explain the output
 to a patient who is not a domain expert.

A number of approaches for XAI emerge from the literature and broadly contain two parts: (1) transparent box design, which includes algorithms such as decision trees, where models can be directly interpreted by users and therefore an explanation of an output results in simply following the decision paths that relate input to output; (2) post hoc interpretability, which provides an explanation to a black-box model using additional methods such as probing, perturbing, or by constructing surrogate models for general ML or AI methods [93, 170].

Utilization of KGs can greatly enhance XAI qualities as KGs are ideal for improving the 673 model's interpretability, explainability, and understandability. Some methods are directly built 674 around KGs and thus take full advantage of them. Examples of those methods may include 675 methods that are using paths [164], predicting links, or performing reasoning [34]. Other methods 676 can be enhanced using the KG (e.g., [129]). Yet the enhancement effect greatly depends on 677 the place where KGs are employed and iteratively applied: pre-model (e.g., KG construction, 678 potentially multi-modal), in-model (e.g., integrating KG with machine learning models), and post-679 model (e.q., reviewing and updating KG by domain experts to be applied in the next iteration 680 to enhance machine learning models and their explanability) [146]. For example in in-model 681 use, a model can be pre-trained using a KG, and an example of a pre-trained language model is 682 SapBERT [111], which utilises synonyms in the UMLS Metathesaurus to further pre-train a BERT 683 language model. This can not only be beneficial for performance [195], but can also potentially 684 enhance post-model explanation since the trained features are aligned with the KG [146]. 685

5.1 What has been done: use cases and recent developments

Explainable AI for Healthcare Practice. The utilization of AI in healthcare practice raises the concern of leaving life-critical decisions to black-box models [146, 170]. For example, in the field of precision medicine which aims at tailoring drug treatments and dosages to each patient, clinicians require more information from a model than a simple binary decision [8]. The interpretability and explainability of AI models is thus an essential characteristic to make outputs understandable and transparent. This would enforce both clinicians' and patients' trust in models by complementing (and not substituting) clinicians' explanations [21, 146, 170].

To illustrate, this direction has been envisioned for several healthcare scenarios. Explainable AI models could support the experts in finding clinical trials that are appropriate based on patient history [170]. Counterintuitive or unreliable predictions that could have serious consequences could be explained, and thus prevented [170, 15, 92]. Some also envision such models to be used to explain and debunk healthcare-related misinformation [146]. As aforementioned, it is noteworthy that different kinds of explanations should be employed depending on the target audience, *e.g.*, scientific explanations for evidence or trace-based explanations for treatment [21].

Explainable AI for Knowledge Discovery. As introduced in Section 4, KGs can support 701 knowledge discovery in life science, including the explainability of the process and the discovered 702 units. In this view, Ritoski and Paulheim [150] explain that ontologies, linked data, and KGs are 703 used in the interpretation step of a data mining process, e.g., for interpreting sequential patterns 704 in patient data [79], or to describe subgroups in a semantic subgroup discovery process [172]. 705 KGs can also serve both as the basis for knowledge discovery processes and the interpretation 706 process. For example, Linked Open Data connecting drugs and adverse reactions can be analyzed 707 with Hidden Conditional Random Fields to predict adverse drug reactions, where the paths from 708 selected drugs to outcomes visually explain the prediction [89]. Similarly, Bresso et al. [13] lever-709

age features extracted from KGs (interpretable features such as paths, neighbors, path patterns) and white box models (*e.g.*, decision trees) to reproduce expert classifications of drugs causing or not specific adverse drug reactions. The rules extracted from the decision trees contain features that provide explanations for the molecular mechanisms behind these adverse reactions according to experts. Sousa *et al.* [162] employ KGs to explain both protein-protein interaction predictions and gene-disease association predictions based on shared semantic aspects.

Explainable AI for KG Construction The final use case considers the situation that XAI 716 is applied to KGs themselves. We discussed the challenge to support human intervention in KG 717 construction in Section 3.4. Recent KG construction gradually relies on data-driven, deep learning 718 based methods to automatically induce knowledge from data. The deep learning models are 719 opaque, and thus the process requires explainability. The resulting KG may not be accountable to 720 be used for downstream applications. Trustworthy KG engineering is proposed in [194] to highlight 721 the importance of embedding explainable AI and human intervention in the KG life cycle. XAI 722 methods have been applied in many NLP related tasks (entity and relation extraction, entity 723 resolution, link prediction, etc.) in KG construction from texts. The XAI methods rely either on 724 feature-based explanations or knowledge-based explanations. While feature-based explanations 725 try to infer explanations from the data or the models' interpretation of the data, knowledge-based 726 explanations aim to interpret the process with rules, reasoning paths, and structured contextual 727 information. Rules and paths have mainly been used for explanation, especially for link prediction, 728 a task comprehensively surveyed in [194]. 729

730 5.2 What are the challenges?

- How to integrate KGs for better XAI, especially with recent deep learning and 731 language model based methods? KG may provide better data provenance for the model 732 output. This can ensure explainability for communicating the model to domain experts in 733 data science applications [8]. In terms of recent generative LLMs, life science KGs, with careful 734 curation based on scientific publications, may help to provide provenance data to the answers 735 generated by LLMs. Studies need to understand to what extent, and how, LLMs can be applied 736 to induce knowledge (e.q., by probing LLMs with biomedical ontologies [67]), which then may737 provide a foundation to create better approaches to integrate KGs with LLMs. Another area 738 is neuro-symbolic methods which may provide models that are inherently more interpretable 739 (see further discussions in Section 6.1). Also, regarding language models (especially LLMs), 740 they are capable of generating fluent texts, which can potentially serve as textual explanation 741 generators from symbolic knowledge for XAI. Meanwhile, a key issue is the hallucination of 742 LLMs, and KGs may support better prompting, fine-tuning and interpretable inference of 743 LLMs for higher decisiveness and trustfulness [137]. 744

How to evaluate XAI methods that involve KG? How to measure the quality of ex-745 planations, to ensure they are corresponding to users? The majority (around 70%) of XAI 746 studies for KG construction do not evaluate the quality of the explanations or only informally 747 visualize or comment on a limited number of cases to show the intuitive outcome [194]. Also, 748 an XAI method needs to consider the target audience, as the explainability is to be finally 749 received by a group of humans [8]. For instance, only a small number of current approaches 750 to XAI for KG construction involve a user study, human evaluation or task-specific met-751 rics [194]. Evaluating the quality of explanations requires some expert evaluation performed 752 as ex-post evaluation, and well-defined metrics are needed for this task. An example is in [59] 753 to use a combination of users' scores for each predicted explanation in a KG link prediction 754 task, where there are multiple possible explanations. More expert validated and automated 755

42:20 Knowledge Graphs for Life Sciences

evaluation methods and associated metrics are required for KG-related XAI.

6 Discussion and Conclusion

In this work, we have summarized the recent developments of KG research in life science on three
 important topics – KG Construction and Management, Life Science Knowledge Discovery, and
 KG for XAI. While each topic has its specific challenges, there are some common challenges and
 trends for the life science KG research in general.

6.1 Overall challenges and trends

Meanwhile, more scalable and efficient knowledge retrieval, query and reasoning systems, including life science KGs and mapping repositories, are still worthy of investigation and development.

Evolution and Quality Assurance of KGs. KGs need to be updated as new data and 765 knowledge are emerging, and the schema and facts can easily become outdated or less useful for 766 existing applications in life sciences. In terms of KG construction, we discussed ontology extension 767 as a use case to address the evolution issue or emergence of new concepts and relations, and also 768 instance matching to extend new instances for the KG. Updating KGs is also a prerequisite for 769 life science knowledge discovery and knowledge discovery methods should be able to support the 770 evolution of KGs with e.g., the capabilities of continuous learning and zero-shot learning. Quality 771 assurance is another issue for KGs, including the tasks of knowledge error detection and correction, 772 knowledge completion, knowledge canonicalization, etc. On the one hand, more effective KG 773 quality assurance methods and systems should be developed, including schema and constraint 774 languages for quality verification and learning-based models for prediction (e.q., [25] combines)775 both for fact correction); on the other hand, knowledge discovery methods should be robust to 776 noisy KGs by investigating e.g., robust KG embeddings and multi-modal representation learning. 777

Heterogeneity in KGs: Multi-domain and Multi-modality. KGs contain heterogeneous 778 information, which brings challenges to their construction, representation, and reasoning. Differ-779 ent schema and data in KGs can have different focuses in their scopes and domains. Integrating 780 data of different domains for building *multi-domain* KGs is difficult with challenges in e.g., ontol-781 ogy and data matching. Besides, recent studies have explored integrating different modalities to 782 construct Multi-modal KGs [27, 125, 179], for instance text [136], images [182], etc. One challenge 783 to address is how to learn effective machine learning models over multi-modal KGs fused from 784 different sources (patients' records, curated knowledge bases, and scientific literature) to support 785 scientific discovery as well as KG construction and management. Another challenge is developing 786 accurate and efficient knowledge representation approaches for texts and images in multi-modal 787 KG construction. For example, careful consideration should be given to when to simply use an 788 annotation property to associate an image with an entity, and when to use a property with specific 789 semantics to connect an image and an entity. 790

Human Interaction and Explainability with KGs. In KG construction, human experts 791 are required for many sub-tasks of KG construction and provide oversight [194]. In life science 792 knowledge discovery, human experts are necessary to finally validate the predicted new knowl-793 edge. The whole process of interacting with KG in life sciences requires explainability, especially 794 when sub-symbolic models (e.g., pre-trained language models) are used. How to generate clear 795 explanations for human interaction and how to evaluate the quality of explanations remains a 796 challenge, as well as how to achieve consensus regarding scientific understanding with automati-797 cally discovered knowledge when organizing knowledge in life science [132]. The recent growth of 798 Neuro-Symbolic methods suggests their support for explainability [91, 92, 153]. A recent survey 799

[92] summarizes XAI in bioinformatics with a chapter on knowledge-based explanations, whereas Karim [91, Chapter 8] provides a neuro-symbolic framework for KG construction and utilisation for medical experts' decision making in the cancer domain. The approach presented in [153] is another recent example of neuro-symbolic integration for image classification with KG-based XAI in the cultural heritage domain.

Personalized and Customized KGs. A key challenge for KG construction is customisation, as 805 we discussed in Section 3, to construct application-oriented KGs, where relevant sub-KGs have to 806 be extracted for large-scale KGs (a.k.a. modualization) and integrated with other knowledge and 807 data from different sources. Besides, many life science KGs are about individuals, e.g., patients in 808 healthcare applications, where Personal Health KG enables the integration of instance-level (or 809 patient-level) information and their computation is required [125]. An example is the Personal 810 Health KG in [22] that supports the dietary recommendation for users, where the construction 811 and population of the KG requires reusing and integrating existing ontologies, dietary guidelines, 812 and time-series patient data. The other examples of KGs integrating patients' EHR data [169, 16] 813 are presented in Section 4.1. In personal KG construction, personal data should be protected. 814 KG scalability should also be considered in order to be used on small devices such as cellphones. 815 This is still a big challenge that has been rarely considered in using KGs in the life sciences. 816

Distributed KGs. The value of healthcare data for improving clinical knowledge and standard of 817 care and the potential of semantic technologies to further enhance it are well recognized. However, 818 a responsible use of healthcare data at the global level (beyond each healthcare provider and 819 even each country) must take into account both legal and ethical issues in data sharing, privacy 820 and security. Distributed knowledge graphs can mitigate these issues, by allowing for access 821 control and privacy protection. Furthermore, distributed knowledge graphs can also address the 822 challenges of scientific data ownership and stewardship by enabling the decentralized publishing 823 of high quality data. Several approaches for federated querying and embedding of knowledge 824 graphs have been proposed in recent years [26, 139, 160], however a wide adoption of semantic 825 technologies in healthcare is still lacking, with a proliferation of terminological standards and a 826 disconnection between data and meaning. 827

Representation Learning with KGs: Symbolic and Sub-symbolic Integration. Across 828 the topics and use cases, we see the importance of transforming symbolic knowledge into sub-829 symbolic representations or combining both representations. The combination of both the neural 830 and the traditional symbolic representation methods leads to a trend in neural-symbolic ap-831 proaches in the field [12]. Recently, Pre-trained and Large Language Models provide new methods 832 to transfer self-supervised learning from a vast amount of corpora to support KG construction, 833 e.g., OntoGPT [18] and OntoLAMA [67]. LLMs are especially good at representing texts of 834 life science publications in sub-symbolic spaces for semantic understanding. KGs may also pro-835 vide a layer of explainability by validating the output of LLMs. A recent survey [137] proposes a 836 roadmap for integrating LLMs and KGs. OntoProtein [195] is a recent example of how to integrate 837 KGs into the process of pre-training LLMs in the bioinformatic domain, thus achieving improved 838 results on protein-related knowledge discovery tasks. Also, geometry-informed representations 839 of more formal KGs, especially in hyperbolic spaces or using complex geometric structures, e.g., 840 [19, 102], can usually represent the structure of the KG with low dimensional vectors. Graph 841 Neural Networks may also support the encoding of KG structures in a more explainable way with 842 logical rules [31]. 843

42:22 Knowledge Graphs for Life Sciences

844 6.2 Conclusion

Knowledge Graphs have become a popular and effective method to represent heterogeneous con-845 cepts, relations, and data in life sciences. They require scalable solutions to represent and reason 846 with heterogeneous data and require constant updates. Throughout this work, we covered the 847 main topics and their corresponding use cases of KGs in multiple life science domains such as pro-848 tein analysis, drug discovery, ecotoxicology, and healthcare, and summarized the corresponding 849 challenges. As new methods in knowledge representation appear, for instance the recent trends 850 of human-in-the-loop, sub-symbolic knowledge representations, pre-trained and large language 85 models, and neuro-symbolic integration, we envisage deeper applications of KGs to life science 852 processes, that support the construction of more applicable KGs and the discovery of more re-853 liable scientific knowledge, with explainability and human interaction better supported. KGs in 854 combination with other modern machine learning and natural language processing techniques will 855 become a foundation for AI for the life sciences. 856

⁸⁵⁷ Appendix A: Terms in Knowledge Graphs and Life Sciences

Below we provide a list of key terms used in this paper, as well as their definitions and explanations. Note we mainly use the original sentences in the sources that are referenced as the
definitions.

Description Logics: a family of knowledge representation languages that can be used to represent knowledge of an application domain. DLs differ from their predecessors, such as semantic networks and frames, in that they are equipped with logic-based semantics, the same semantics as that of classical first-order logic. Most ontologies are implemented in OWL, whose semantics are given by the Description Logic *SROIQ*. [6]

TBox and **ABox**: the two components of domain knowledge in Description Logics, *i.e.*, a terminological part called the TBox and an assertional part called the ABox, with the combination of a TBox and an ABox being called a knowledge base (KB). The TBox represents knowledge about the structure of the domain (similar to a database schema), while the ABox represents knowledge about a concrete situation (similar to a database instance). [6]

Semantic Networks: a graph structure for representing knowledge in patterns of interconnected nodes and arcs [163]. We use the term to denote a graph of concepts and relations without
formal semantics.

Gene Ontology: The Gene Ontology (GO) knowledgebase provides a comprehensive, structured, computer-accessible representation of gene function, for genes from any cellular organism or virus [5, 29].

SNOMED-CT: Systematized Nomenclature of Medicine Clinical Terms (SNOMED CT) is a structured clinical vocabulary. It has a general and comprehensive coverage of clinical terms to support electronic healthcare systems and clinical applications. [40, 28]

UMLS (UMLS Metahesaurus and UMLS Semantic Networks): Unified Medical Lan guage System (UMLS) is a repository of biomedical vocabularies developed by the US National
 Library of Medicine. The UMLS is composed of three "knowledge sources", a Metathesaurus,
 a semantic network, and a lexicon. The UMLS Metathesaurus is a comprehensive effort for in tegrating biomedical ontologies through mappings. The UMLS Semantic Networks define the
 types or categories, or Semantic Types, of all Metathesaurus concepts and their relationships, or
 Semantic Relations. [10, 28]

ChEBI: Chemical Entities of Biological Interest (ChEBI) is a database and ontology containing information about chemical entities of biological interest. [64]

Symbolic vs. subsymbolic representations: Rooted in cognitive science, symbolic sys-

Chen, Dong, Hastings, Jiménez-Ruiz, Lopez, Monnin, Pesquita, Škoda, Tamma

tems of human cognition are related to the representation and manipulation of symbols; subsymbolic or connectionist systems are most generally associated with the metaphor of a neuron, *e.g.*, perceptrons as an early system [97]. In terms of AI, symbolic systems contain logic-based and knowledge representations, while subsymbolic systems typically contain neural networks and deep learning based methods [43]. Neural language models and pre-trained language models [88] are also classified under subsymbolic systems.

Pre-trained and Large Language Models: Neural language modelling is the task of 896 using neural network approaches to predict words from prior their contexts in a sequence. Pre-897 training is the process of learning some sort of representation (usually neural embedding based) 898 of meaning for words or sentences by processing very large amounts of text (or other data in a 899 sequence form, e.g., proteins and KG facts). This results in pre-trained language models. The 900 dominating architecture for neural language modeling is Transformer-based models, including 901 BERT, its domain specific versions, and later large variants, like the GPT series. The pre-trained 902 language models of very large sizes are recently coined Large Language Models (LLMs). [88] 903

Neuro-symbolic representations: refers to the integration of neural networks and symbolic
 representations to design AI models that base their prediction on both data and knowledge. [43]

906 Appendix B: Authors' Contributions

All authors participated in the planning and discussions of this work. JH and HD finished the abstract and "Introduction". VT, JC and EJR contributed to "Knowledge Graphs in the Life Sciences". VT contributed to the main part of "Knowledge Graph Construction and Management", with contributions of use cases from JC, HD, PM, EJR, and JH. VL and JC contributed to "Life Science Knowledge Discovery". PM, PS, HD, and CP contributed to "Knowledge Graphs for Explainable AI". HD, JC, and CP contributed to "Discussion and Conclusion" based on discussions with other team members. All authors contributed to the final revision of this paper.

Acknowledgements We would like to thank Uli Sattler (University of Manchester) for proposing
 the topic of this paper and Terry Payne (University of Liverpool) for the useful comments on a
 previous draft. We would also like to thank the TGDK editors in chief for organizing this inaugural
 issue.

— References

- 1Daniel J Abadi, Adam Marcus, Samuel R Madden, and Kate Hollenbach. SW-Store: a vertically partitioned DBMS for Semantic Web data management. *The VLDB Journal*, 18:385–406, 2009.
- 2Emily Alsentzer, Michelle M. Li, Shilpa N. Kobren, Undiagnosed Diseases Network, Isaac S. Kohane, and Marinka Zitnik. Deep learning for diagnosing patients with rare genetic diseases. medRxiv, 2022.
- 3Robert Arp, Barry Smith, and Andrew D. Spear. Building Ontologies With Basic Formal Ontology. The MIT Press, 08 2015.
- 4M Ashburner, C A Ball, J A Blake, D Botstein, H Butler, J M Cherry, A P Davis, K Dolinski, S S Dwight, J T Eppig, M A Harris, D P Hill, L Issel-Tarver, A Kasarskis, S Lewis, J C Matese, J E Richardson, M Ringwald, G M Rubin, and G Sherlock. Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. Nat Genet, 25(1):25–29, May 2000.
- 5 Michael Ashburner, Catherine A Ball, Judith A Blake, David Botstein, Heather Butler, J Michael

Cherry, Allan P Davis, Kara Dolinski, Selina S Dwight, Janan T Eppig, et al. Gene ontology: tool for the unification of biology. *Nature genetics*, 25(1):25–29, 2000.

- 6Franz Baader, Ian Horrocks, Carsten Lutz, and Uli Sattler. An Introduction to Description Logic. Cambridge University Press, Cambridge, 2017.
- 7Nguyen Bach and Sameer Badaskar. A review of relation extraction. *Literature review for Language* and Statistics II, 2:1–15, 2007.
- 8Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bennetot, Siham Tabik, Alberto Barbado, Salvador Garcia, Sergio Gil-Lopez, Daniel Molina, Richard Benjamins, Raja Chatila, and Francisco Herrera. Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. Information Fusion, 58:82–115, 2020.
- 9 Maciej Besta, Robert Gerstenberger, Emanuel Peter, Marc Fischer, Michał Podstawski, Claude

Barthels, Gustavo Alonso, and Torsten Hoefler. Demystifying graph databases: Analysis and taxonomy of data organization, system designs, and graph queries. *ACM Comput. Surv.*, 2023.

- 10Olivier Bodenreider. The Unified Medical Language System (UMLS): integrating biomedical terminology. Nucleic Acids Res., 32(Database-Issue):267–270, 2004.
- 11Nadav Brandes, Dan Ofer, Yam Peleg, Nadav Rappoport, and Michal Linial. ProteinBERT: a universal deep-learning model of protein sequence and function. *Bioinformatics*, 38(8):2102–2110, 02 2022.
- 12 Anna Breit, Laura Waltersdorfer, Fajar J. Ekaputra, Marta Sabou, Andreas Ekelhart, Andreea Iana, Heiko Paulheim, Jan Portisch, Artem Revenko, Frank van Harmelen, and Annette ten Teije. Combining machine learning and semantic web: A systematic mapping study. ACM Computing Surveys, 2023.
- 13Emmanuel Bresso, Pierre Monnin, Cédric Bousquet, François-Élie Calvier, Ndeye Coumba Ndiaye, Nadine Petitpain, Malika Smaïl-Tabbone, and Adrien Coulet. Investigating ADR mechanisms with explainable AI: a feasibility study with knowledge graph mining. *BMC Medical Informatics Decis. Mak.*, 21(1):171, 2021.
- 14Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020.
- 15Rich Caruana, Yin Lou, Johannes Gehrke, Paul Koch, Marc Sturm, and Noemie Elhadad. Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission. In Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pages 1721–1730, 2015.
- 16 Ricardo MS Carvalho, Daniela Oliveira, and Catia Pesquita. Knowledge Graph Embeddings for ICU readmission prediction. BMC Medical Informatics and Decision Making, 23(1):12, 2023.
- 17 Javier Castell-Díaz, Jose Antonio Miñarro-Giménez, and Catalina Martínez-Costa. Supporting SNOMED CT postcoordination with knowledge graph embeddings. Journal of Biomedical Informatics, 139:104297, 2023.
- 18 J. Harry Caufield, Harshad Hegde, Vincent Emonet, Nomi L. Harris, Marcin P. Joachimiak, Nicolas Matentzoglu, HyeongSik Kim, Sierra A. T. Moxon, Justin T. Reese, Melissa A. Haendel, Peter N. Robinson, and Christopher J. Mungall. Structured prompt interrogation and recursive extraction of semantics (SPIRES): A method for populating knowledge bases using zero-shot learning, April 2023.
- 19 Ines Chami, Adva Wolf, Da-Cheng Juan, Frederic Sala, Sujith Ravi, and Christopher Ré. Lowdimensional hyperbolic knowledge graph embeddings. arXiv preprint arXiv:2005.00545, 2020.
- 20Payal Chandak, Kexin Huang, and Marinka Zitnik. Building a knowledge graph to enable precision medicine. *bioRxiv*, 2022.

- 21Shruthi Chari, Oshani Seneviratne, Daniel M. Gruen, Morgan A. Foreman, Amar K. Das, and Deborah L. McGuinness. Explanation ontology: A model of explanations for user-centered AI. In Jeff Z. Pan, Valentina A. M. Tamma, Claudia d'Amato, Krzysztof Janowicz, Bo Fu, Axel Polleres, Oshani Seneviratne, and Lalana Kagal, editors, The Semantic Web - ISWC 2020 - 19th International Semantic Web Conference, Athens, Greece, November 2-6, 2020, Proceedings, Part II, volume 12507 of Lecture Notes in Computer Science, pages 228–243, 2020.
- 22 Ching-Hua Chen, Daniel Gruen, Jonathan Harris, James Hendler, Deborah L McGuinness, Marco Monti, Nidhi Rastogi, Oshani Seneviratne, and Mohammed J Zaki. Semantic technologies for clinically relevant personal health applications. In *Personal Health Informatics: Patient Participation in Precision Health*, pages 199–220. Cham, 2022.
- 23 Jiaoyan Chen, Yuxia Geng, Zhuo Chen, Jeff Z Pan, Yuan He, Wen Zhang, Ian Horrocks, and Huajun Chen. Zero-Shot and Few-Shot Learning With Knowledge Graphs: A Comprehensive Survey. *Proceedings of the IEEE*, 2023.
- 24 Jiaoyan Chen, Yuan He, Yuxia Geng, Ernesto Jiménez-Ruiz, Hang Dong, and Ian Horrocks. Contextual semantic embeddings for ontology subsumption prediction. World Wide Web, pages 1– 23, 2023.
- 25 Jiaoyan Chen, Ernesto Jiménez-Ruiz, Ian Horrocks, Xi Chen, and Erik Bryhn Myklebust. An assertion and alignment correction framework for large scale knowledge bases. *Semantic Web*, 14(1):29–53, 2023.
- 26 Mingyang Chen, Wen Zhang, Zonggang Yuan, Yantao Jia, and Huajun Chen. Federated knowledge graph completion via embedding-contrastive learning. *Knowledge-Based Systems*, 252:109459, 2022.
- 27 Yong Chen, Xinkai Ge, Shengli Yang, Linmei Hu, Jie Li, and Jinwen Zhang. A survey on multimodal knowledge graphs: Construction, completion and applications. *Mathematics*, 11(8):1815, 2023.
- 28E. Coiera. Guide to Health Informatics, chapter Chapter 23 Healthcare terminologies and classification systems, pages 381–399. CRC Press, Taylor & Francis Group, Boca Raton, 2015.
- 29 The Gene Ontology Consortium, Suzi A Aleksander, James Balhoff, Seth Carbon, J Michael Cherry, Harold J Drabkin, Dustin Ebert, Marc Feuermann, Pascale Gaudet, Nomi L Harris, David P Hill, Raymond Lee, Huaiyu Mi, Sierra Moxon, Christopher J Mungall, Anushya Muruganugan, Tremayne Mushayahama, Paul W Sternberg, Paul D Thomas, Kimberly Van Auken, Jolene Ramsey, Deborah A Siegele, Rex L Chisholm, Petra Fey, Maria Cristina Aspromonte, Maria Victoria Nugnes, Federica Quaglia, Silvio Tosatto, Michelle Giglio, Suvarna Nadendla, Giulia Antonazzo, Helen Attrill, Gil dos Santos, Steven Marygold, Victor Strelets, Christopher J Tabone, Jim Thurmond, Pinglei Zhou, Saadullah H Ahmed, Praoparn Asanitthong, Diana Luna Buitrago, Meltem N Erdol, Matthew C Gage, Mohamed Ali Kadhum, Kan Yan Chloe Li, Miao Long, Aleksandra Michalak, Angeline

34Ivan Donadello, Mauro Dragoni, and Claudio Eccher. Persuasive explanation of reasoning inferences on dietary data. In PROFILES/SE-MEX@JCHUC 2010

35Hang Dong, Jiaoyan Chen, Yuan He, and Ian Horrocks. Ontology enrichment from texts: A biomedical dataset for concept discovery and placement. In Proceedings of the 32nd ACM International Conference on Information & Knowledge Management, 2023.

MEX@ISWC, 2019.

- 36Hang Dong, Jiaoyan Chen, Yuan He, Yinan Liu, and Ian Horrocks. Reveal the unknown: Out-ofknowledge-base mention discovery with entity linking. In Proceedings of the 32nd ACM International Conference on Information & Knowledge Management, 2023.
- 37 Hang Dong, Matúš Falis, William Whiteley, Beatrice Alex, Joshua Matterson, Shaoxiong Ji, Jiaoyan Chen, and Honghan Wu. Automated clinical coding: what, why, and where we are? NPJ digital medicine, 5(1):159, 2022.
- 38Hang Dong, Víctor Suárez-Paniagua, Huayu Zhang, Minhong Wang, Arlene Casey, Emma Davidson, Jiaoyan Chen, Beatrice Alex, William Whiteley, and Honghan Wu. Ontology-driven and weakly supervised rare disease identification from clinical notes. BMC Medical Informatics and Decision Making, 23(1):1–17, 2023.
- 39Kevin Donnelly et al. SNOMED-CT: The advanced terminology and coding system for ehealth. In Medical and Care Computerics 3, volume 121 of Studies in health technology and informatics, pages 279–290. IOS Press, 2006.
- 40Kevin Donnelly et al. SNOMED-CT: The advanced terminology and coding system for ehealth. In Medical and Care Computers 3, volume 121 of Studies in health technology and informatics, pages 279–290. IOS Press, 2006.
- 41Zlatan Dragisic, Valentina Ivanova, Huanyu Li, and Patrick Lambrix. Experiences from the anatomy track in the ontology alignment evaluation initiative. J. Biomed. Semant., 8(1):56:1– 56:28, 2017.
- 42Juan M Durán. Dissecting scientific explanation in ai (sxai): A case for medicine and healthcare. *Artificial Intelligence*, 297:103498, 2021.
- 43Artur S d'Avila Garcez, Luís C Lamb, and Dov M Gabbay. Neural-symbolic learning systems, pages 35–54. Springer, 2009.
- 44 Joakim Edin, Alexander Junge, Jakob D. Havtorn, Lasse Borgholt, Maria Maistro, Tuukka Ruotsalo, and Lars Maaløe. Automated medical coding on MIMIC-III and MIMIC-IV: A critical review and replicability study. CoRR, abs/2304.10909, 2023.
- 45Lisa Ehrlinger and Wolfram Wöß. Towards a definition of knowledge graphs. In International Conference on Semantic Systems, 2016.
- 46Daniel Faria, Ernesto Jiménez-Ruiz, Catia Pesquita, Emanuel Santos, and Francisco M. Couto. Towards Annotating Potential Incoherences in BioPortal Mappings. In Peter Mika, Tania Tudorache, Abraham Bernstein, Chris Welty, Craig A. Knoblock, Denny Vrandecic, Paul Groth, Natasha F. Noy, Krzysztof Janowicz, and Carole A. Goble, editors, The Semantic Web -ISWC 2014 - 13th International Semantic Web

James Seager, Laurel Cooper, Justin Elser, Pankaj Jaiswal, Parul Gupta, Pankaj Jaiswal, Sushma Naithani, Manuel Lera-Ramirez, Kim Rutherford, Valerie Wood, Jeffrey L De Pons, Melinda R Dwinell, G Thomas Hayman, Mary L Kaldunski, Anne E Kwitek, Stanley J F Laulederkind, Marek A Tutaj, Mahima Vedi, Shur-Jen Wang, Peter D'Eustachio, Lucila Aimo, Kristian Axelsen, Alan Bridge, Nevila Hyka-Nouspikel, Anne Morgat, Suzi A Aleksander, J Michael Cherry, Stacia R Engel, Kalpana Karra, Stuart R Miyasato, Robert S Nash, Marek S Skrzypek, Shuai Weng, Edith D Wong, Erika Bakker, Tanya Z Berardini, Leonore Reiser, Andrea Auchincloss, Kristian Axelsen, Ghislaine Argoud-Puy, Marie-Claude Blatter, Emmanuel Boutet, Lionel Breuza, Alan Bridge, Cristina Casals-Casas, Elisabeth Coudert, Anne Estreicher, Maria Livia Famiglietti, Marc Feuermann, Arnaud Gos, Nadine Gruaz-Gumowski, Chantal Hulo, Nevila Hyka-Nouspikel, Florence Jungo, Philippe Le Mercier, Damien Lieberherr, Patrick Masson, Anne Morgat, Ivo Pedruzzi, Lucille Pourcel, Sylvain Poux, Catherine Rivoire, Shyamala Sundaram, Alex Bateman, Emily Bowler-Barnett, Hema Bye-A-Jee, Paul Denny, Alexandr Ignatchenko, Rizwan Ishtiaq, Antonia Lock, Yvonne Lussi, Michele Magrane, Maria J Martin, Sandra Orchard, Pedro Raposo, Elena Speretta, Nidhi Tyagi, Kate Warner, Rossana Zaru, Alexander D Diehl, Raymond Lee, Juancarlos Chan, Stavros Diamantakis, Daniela Raciti, Magdalena Zarowiecki, Malcolm Fisher, Christina James-Zorn, Virgilio Ponferrada, Aaron Zorn, Sridhar Ramachandran, Levla Ruzicka, and Monte Westerfield. The Gene Ontology knowledgebase in 2023. Genetics, 224(1):iyad031, 03 2023.

Pesala, Armalya Pritazahra, Shirin C C Saver-

imuttu, Renzhi Su, Kate E Thurlow, Ruth C

Lovering, Colin Logie, Snezhana Oliferenko, Ju-

dith Blake, Karen Christie, Lori Corbani, Mary E

Dolan, Harold J Drabkin, David P Hill, Li Ni,

Dmitry Sitnikov, Cynthia Smith, Alayne Cuzick,

- 30 The UniProt Consortium. UniProt: a worldwide hub of protein knowledge. Nucleic Acids Research, 47(D1):D506-D515, 11 2018.
- 31 David Jaime Tena Cucala, Bernardo Cuenca Grau, Egor V Kostylev, and Boris Motik. Explainable GNN-Based Models over Knowledge Graphs. In International Conference on Learning Representations, 2021.
- 32Shumin Deng, Chengming Wang, Zhoubo Li, Ningyu Zhang, Zelin Dai, Hehong Chen, Feiyu Xiong, Ming Yan, Qiang Chen, Mosha Chen, et al. Construction and applications of billionscale pre-trained multimodal business knowledge graph. In 2023 IEEE 39th International Conference on Data Engineering (ICDE), pages 2988– 3002. IEEE, 2023.
- 33 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), 2019.

Conference. Proceedings, Part II, volume 8797 of Lecture Notes in Computer Science, pages 17–32, 2014.

- 47 Mariano Fernandez, Asuncion Gomez-Perez, and Natalia Juristo. Methontology: from ontological art towards ontological engineering. In Proc. of the AAAI97 Spring Symposium Series on Ontological Engineering, pages 33–40. Stanford, USA, 1997.
- 48 Giorgos Flouris, Dimitris Manakanatas, Haridimos Kondylakis, Dimitris Plexousakis, and Grigoris Antoniou. Ontology change: classification and survey. *Knowl. Eng. Rev.*, 23(2):117–152, 2008.
- 49 Nadime Francis, Alastair Green, Paolo Guagliardo, Leonid Libkin, Tobias Lindaaker, Victor Marsault, Stefan Plantikow, Mats Rydberg, Petra Selmer, and Andrés Taylor. Cypher: An Evolving Query Language for Property Graphs. In Proceedings of the 2018 International Conference on Management of Data, SIGMOD '18, page 1433–1445, 2018.
- 50 Katrin Fundel, Robert Küffner, and Ralf Zimmer. Relex—relation extraction using dependency parse trees. *Bioinformatics*, 23(3):365–371, 2007.
- 51 Thomas Gaudelet, Ben Day, Arian R Jamasb, Jyothish Soman, Cristian Regep, Gertrude Liu, Jeremy B R Hayter, Richard Vickers, Charles Roberts, Jian Tang, David Roblin, Tom L Blundell, Michael M Bronstein, and Jake P Taylor-King. Utilizing graph machine learning within drug discovery and development. *Briefings in Bioinformatics*, 22(6):bbab159, 05 2021.
- 52 Anna Gaulton, Louisa J. Bellis, A. Patricia Bento, Jon Chambers, Mark Davies, Anne Hersey, Yvonne Light, Shaun McGlinchey, David Michalovich, Bissan Al-Lazikani, and John P. Overington. Chembl: A large-scale bioactivity database for drug discovery. Nucleic acids research, 40(D1):D1100–D1107, 2012.
- 53David Geleta, Andriy Nikolov, Gavin Edwards, Anna Gogleva, Richard Jackson, Erik Jansson, Andrej Lamov, Sebastian Nilsson, Marina Pettersson, Vladimir Poroshin, et al. Biological insights knowledge graph: an integrated knowledge graph to support drug development. *Biorxiv*, pages 2021–10, 2021.
- 54Aryo Pradipta Gema, Dominik Grabarczyk, Wolf De Wulf, Piyush Borole, Javier Antonio Alfaro, Pasquale Minervini, Antonio Vergari, and Ajitha Rajan. Knowledge Graph Embeddings in the Biomedical Domain: Are They Useful? A Look at Link Prediction, Rule Learning, and Downstream Polypharmacy Tasks, 2023. arXiv:2305. 19979.
- 55 Martin Glauer, Adel Memariani, Fabian Neuhaus, Till Mossakowski, and Janna Hastings. Interpretable Ontology Extension in Chemistry. Semantic Web Journal, 2022.
- 56 Anika Groß, Cédric Pruski, and Erhard Rahm. Evolution of biomedical ontologies and mappings: Overview of recent approaches. *Computational* and Structural Biotechnology Journal, 14:333–340, 2016.
- 57 Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Fosca Giannotti, and Dino Pe-

dreschi. A survey of methods for explaining black box models. *ACM Comput. Surv.*, 51(5), aug 2018.

- 58Ricardo Guimarães and Ana Ozaki. Reasoning in Knowledge Graphs. In International Research School in Artificial Intelligence in Bergen (AIB 2022), volume 99 of Open Access Series in Informatics (OASIcs), pages 2:1–2:31, 2022.
- 59Nicholas Halliwell, Fabien Gandon, and Freddy Lécué. User scored evaluation of non-unique explanations for relational graph convolutional network link prediction on knowledge graphs. In Anna Lisa Gentile and Rafael Gonçalves, editors, K-CAP '21: Knowledge Capture Conference, Virtual Event, USA, December 2-3, 2021, pages 57– 64, 2021.
- 60I. Harrow, R. Balakrishnan, E. Jimenez-Ruiz, S. Jupp, J. Lomax, J. Reed, M. Romacker, C. Senger, A. Splendiani, J. Wilson, and P. Woollard. Ontology mapping for semantically enabled applications. *Drug Discovery Today*, May 2019.
- 61 Ian Harrow, Ernesto Jiménez-Ruiz, Andrea Splendiani, Martin Romacker, Peter Woollard, Scott Markel, Yasmin Alam-Faruque, Martin Koch, James Malone, and Arild Waaler. Matching disease and phenotype ontologies in the ontology alignment evaluation initiative. J. Biomed. Semant., 8(1):55:1–55:13, 2017.
- **62**J. Hastings. *AI for Scientific Discovery*. AI for everything series. CRC Press, Milton, 2023.
- 63 Janna Hastings. Primer on Ontologies. In Christophe Dessimoz and Nives Škunca, editors, *The Gene Ontology Handbook*, volume 1446, pages 3–13. Humana Press, SpringerOpen, New York, New York, NY, 2017.
- 64Janna Hastings, Gareth Owen, Adriano Dekker, Marcus Ennis, Namrata Kale, Venkatesh Muthukrishnan, Steve Turner, Neil Swainston, Pedro Mendes, and Christoph Steinbeck. Chebi in 2016: Improved services and an expanding collection of metabolites. Nucleic acids research, 44(D1):D1214–D1219, 2016.
- 65 Yuan He, Jiaoyan Chen, Hang Dong, Ian Horrocks, Carlo Allocca, Taehun Kim, and Brahmananda Sapkota. DeepOnto: A Python package for ontology engineering with deep learning. arXiv preprint arXiv:2307.03067, 2023.
- 66 Yuan He, Jiaoyan Chen, Hang Dong, Ernesto Jiménez-Ruiz, Ali Hadian, and Ian Horrocks. Machine learning-friendly biomedical datasets for equivalence and subsumption ontology matching. In Ulrike Sattler, Aidan Hogan, C. Maria Keet, Valentina Presutti, João Paulo A. Almeida, Hideaki Takeda, Pierre Monnin, Giuseppe Pirrò, and Claudia d'Amato, editors, The Semantic Web - ISWC 2022 - 21st International Semantic Web Conference, Proceedings, volume 13489 of Lecture Notes in Computer Science, pages 575-591, 2022.
- 67 Yuan He, Jiaoyan Chen, Ernesto Jiménez-Ruiz, Hang Dong, and Ian Horrocks. Language model analysis for ontology subsumption inference. Findings of the Association for Computational Linguistics: ACL 2023, pages 3439–3453, 2023.
- 68Nicolas Heist and Heiko Paulheim. NASTyLinker: NIL-Aware Scalable Transformer-Based Entity Linker. In Catia Pesquita, Ernesto Jiménez-Ruiz,

Jamie P. McCusker, Daniel Faria, Mauro Dragoni, Anastasia Dimou, Raphaël Troncy, and Sven Hertling, editors, *The Semantic Web - 20th International Conference, ESWC 2023, Hersonissos, Crete, Greece, May 28 - June 1, 2023, Proceedings*, volume 13870 of *Lecture Notes in Computer Science*, pages 174–191. Springer, 2023.

- 69Daniel Scott Himmelstein, Antoine Lizee, Christine Hessler, Leo Brueggeman, Sabrina L Chen, Dexter Hadley, Ari Green, Pouya Khankhanian, and Sergio E Baranzini. Systematic integration of biomedical knowledge prioritizes drugs for repurposing. *Elife*, 6:e26726, 2017.
- 70 Robert Hoehndorf, Paul N. Schofield, and Georgios V. Gkoutos. The role of ontologies in biological and biomedical research: a functional perspective. *Briefings in Bioinformatics*, 16(6):1069–1080, 04 2015.
- 71Aidan Hogan, Eva Blomqvist, Michael Cochez, Claudia d'Amato, Gerard de Melo, Claudio Gutiérrez, Sabrina Kirrane, José Emilio Labra Gayo, Roberto Navigli, Sebastian Neumaier, Axel-Cyrille Ngonga Ngomo, Axel Polleres, Sabbir M. Rashid, Anisa Rula, Lukas Schmelzeisen, Juan F. Sequeda, Steffen Staab, and Antoine Zimmermann. Knowledge Graphs. Number 22 in Synthesis Lectures on Data, Semantics, and Knowledge. Morgan & Claypool Publishers, 2022. URL: https://kgbook.org/.
- 72 Chao-Wei Huang, Shang-Chi Tsai, and Yun-Nung Chen. PLM-ICD: Automatic ICD coding with pretrained language models. In *Proceedings of the 4th Clinical Natural Language Processing Workshop*, pages 10–20, 2022.
- 73 Chung-Chi Huang and Zhiyong Lu. Community challenges in biomedical text mining over 10 years: success, failure and the future. *Briefings in bioinformatics*, 17(1):132–144, 2016.
- 74Kexin Huang, Payal Chandak, Qianwen Wang, Shreyas Havaldar, Akhil Vaid, Jure Leskovec, Girish Nadkarni, Benjamin Glicksberg, Nils Gehlenborg, and Marinka Zitnik. Zero-shot prediction of therapeutic use with geometric deep learning and clinician centered design. medRxiv, 2023.
- 75Kexin Huang, Tianfan Fu, Wenhao Gao, Yue Zhao, Yusuf Roohani, Jure Leskovec, Connor W. Coley, Cao Xiao, Jimeng Sun, and Marinka Zitnik. Artificial intelligence foundation for therapeutic science. *Nature chemical biology*, 2022.
- 76 Pere-Lluís Huguet Cabot and Roberto Navigli. REBEL: Relation extraction by end-to-end language generation. In *Findings of the Association* for Computational Linguistics: EMNLP 2021, pages 2370–2381, Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.
- 77L Hunter, Z Lu, and J Firby. Wab jr, hl johnson, pv ogren, and kb cohen, "opendmap: An open source, ontologydriven concept analysis engine, with applications to capturing knowledge regarding protein transport, protein interactions and cell-typespecific gene expression,". BMC Bioinformatics, 9(78), 2008.
- 78 Anastasiia Iurshina, Jiaxin Pan, Rafika Boutalbi, and Steffen Staab. NILK: Entity Linking Dataset Targeting NIL-Linking Cases. In Proceedings of the

31st ACM International Conference on Information & Knowledge Management, CIKM '22, page 4069–4073, New York, NY, USA, 2022. Association for Computing Machinery.

- 79Nicolas Jay and Mathieu d'Aquin. Linked data and online classifications to organise mined patterns in patient data. In AMIA 2013, American Medical Informatics Association Annual Symposium, Washington, DC, USA, November 16-20, 2013, 2013.
- **80**Shaoxiong Ji, Matti Hölttä, and Pekka Marttinen. Does the magic of BERT apply to medical code assignment? A quantitative study. *Computers in biology and medicine*, 139:104998, 2021.
- 81Shaoxiong Ji, Wei Sun, Hang Dong, Honghan Wu, and Pekka Marttinen. A unified review of deep learning for automated medical coding. arXiv preprint arXiv:2201.02797, 2022.
- 82Ernesto Jiménez-Ruiz and Bernardo Cuenca Grau. LogMap: Logic-based and scalable ontology matching. In The Semantic Web–ISWC 2011: 10th International Semantic Web Conference, Bonn, Germany, October 23-27, 2011, Proceedings, Part I 10, pages 273–288, 2011.
- 83Ernesto Jiménez-Ruiz, Bernardo Cuenca Grau, Ian Horrocks, and Rafael Berlanga Llavori. Logicbased assessment of the compatibility of UMLS ontology sources. J. Biomed. Semant., 2(S-1):S2, 2011. URL: http://www.jbiomedsem.com/ content/2/S1/S2.
- 84Ernesto Jiménez-Ruiz, Bernardo Cuenca Grau, Ulrike Sattler, Thomas Schneider, and Rafael Berlanga Llavori. Safe and Economic Re-Use of Ontologies: A Logic-Based Methodology and Tool Support. In The Semantic Web: Research and Applications, 5th European Semantic Web Conference, Proceedings, pages 185–199, 2008.
- 85Ernesto Jiménez-Ruiz, Christian Meilicke, Bernardo Cuenca Grau, and Ian Horrocks. Evaluating mapping repair systems with large biomedical ontologies. In Thomas Eiter, Birte Glimm, Yevgeny Kazakov, and Markus Krötzsch, editors, Informal Proceedings of the 26th International Workshop on Description Logics, Ulm, Germany, July 23 - 26, 2013, volume 1014 of CEUR Workshop Proceedings, pages 246–257. CEUR-WS.org, 2013.
- **86**Simon Jupp, Thomas Liener, Sirarat Sarntivijai, Olga Vrousgou, Tony Burdett, and Helen E. Parkinson. OxO - A Gravy of Ontology Mapping Extracts. In Matthew Horridge, Phillip Lord, and Jennifer D. Warrender, editors, *Proceedings of the* 8th International Conference on Biomedical Ontology (ICBO 2017), volume 2137 of CEUR Workshop Proceedings, 2017.
- 87Simon Jupp, James Malone, Jerven T. Bolleman, Marco Brandizi, Mark Davies, Leyla J. García, Anna Gaulton, Sebastien Gehant, Camille Laibe, Nicole Redaschi, Sarala M. Wimalaratne, Maria Jesus Martin, Nicolas Le Novère, Helen E. Parkinson, Ewan Birney, and Andrew M. Jenkinson. The EBI RDF platform: linked open data for the life sciences. *Bioinform.*, 30(9):1338–1339, 2014.

88 Daniel Jurafsky and James H. Martin. Speech and Language Processing (3rd Edition). Online, 2023.

- 89 Maulik R. Kamdar and Mark A. Musen. PhLeGrA: Graph Analytics in Pharmacology over the Web of Life Sciences Linked Open Data. In Rick Barrett, Rick Cummings, Eugene Agichtein, and Evgeniy Gabrilovich, editors, *Proceedings of the 26th In*ternational Conference on World Wide Web, pages 321–329. ACM, 2017.
- 90 Maulik R. Kamdar and Mark A. Musen. An empirical meta-analysis of the life sciences linked open data on the web. *Scientific Data*, 8, 2020.
- 91 Md. Rezaul Karim. Interpreting black-box machine learning models with decision rules and knowledge graph reasoning. Dissertation, RWTH Aachen University, Aachen, 2022. Veröffentlicht auf dem Publikationsserver der RWTH Aachen University; Dissertation, RWTH Aachen University, 2022. doi: 10.18154/RWTH-2022-07610.
- 92Md. Rezaul Karim, Tanhim Islam, Oya Beyan, Christoph Lange, Michael Cochez, Dietrich Rebholz-Schuhmann, and Stefan Decker. Explainable AI for bioinformatics: Methods, tools, and applications. *CoRR*, abs/2212.13261, 2022.
- 93Md. Rezaul Karim, Tanhim Islam, Oya Beyan, Christoph Lange, Michael Cochez, Dietrich Rebholz-Schuhmann, and Stefan Decker. Explainable AI for Bioinformatics: Methods, Tools, and Applications, 2023. arXiv:2212.13261.
- 94 Nora Kassner, Fabio Petroni, Mikhail Plekhanov, Sebastian Riedel, and Nicola Cancedda. EDIN: An end-to-end benchmark and pipeline for unknown entity discovery and indexing. In *Proceedings of* the 2022 Conference on Empirical Methods in Natural Language Processing, pages 8659–8673, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics.
- **95**Rohit J Kate. Automatic full conversion of clinical terms into SNOMED CT concepts. *Journal of Biomedical Informatics*, 111:103585, 2020.
- 96 Mayank Kejriwal. Domain-Specific Knowledge Graph Construction. Springer Publishing Company, Incorporated, 1st edition, 2019.
- 97 Troy D. Kelley. Symbolic and sub-symbolic representations in computational models of human cognition: What can be learned from biology? *Theory* & *Psychology*, 13(6):847–860, 2003.
- 98Elisa F. Kendall and Deborah L. McGuinness. Ontology Engineering. Synthesis Lectures on the Semantic Web: Theory and Technology. Springer, Cham, Switzerland, 2019.
- 99Sebastian Köhler, Michael Gargano, Nicolas Matentzoglu, Leigh C Carmody, David Lewis-Smith, Nicole A Vasilevsky, Daniel Danis, Ganna Balagura, Gareth Baynam, Amy M Brower, et al. The human phenotype ontology in 2021. Nucleic acids research, 49(D1):D1207–D1217, 2021.
- 100 Martin Krallinger, Florian Leitner, Carlos Rodriguez-Penagos, and Alfonso Valencia. Overview of the protein-protein interaction annotation extraction task of BioCreative II. Genome biology, 9:1–19, 2008.
- 101 Michael Kuhn, Damian Milosz Szklarczyk, Sune Pletscher-Frankild, Thomas H Blicher, Christian von Mering, Lars J Jensen, and Peer Bork. Stitch 4: integration of protein-chemical interactions with

user data. Nucleic Acids Research, 42(D1):D401–D407, 2013.

- 102Maxat Kulmanov and Robert Hoehndorf. Deep-GOZero: improving protein function prediction from sequence and zero-shot learning based on ontology axioms. *Bioinformatics*, 38(Supplement_1):i238–i245, 2022.
- 103Hoang Thanh Lam, Marco Luca Sbodio, Marcos Martinez Gallindo, Mykhaylo Zayats, Raul Fernandez-Diaz, Victor Valls, Gabriele Picco, Cesar Berrospi Ramis, and Vanessa Lopez. Otter-Knowledge: benchmarks of multimodal knowledge graph representation learning from different sources for drug discovery. CoRR, abs/2306.12802, 2023.
- 104Robert Leaman, Ritu Khare, and Zhiyong Lu. Challenges in clinical natural language processing for automated disorder normalization. *Journal of biomedical informatics*, 57:28–37, 2015.
- 105 Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240, 09 2019.
- 106 Ulf Leser and Jörg Hakenberg. What makes a gene name? Named entity recognition in the biomedical literature. Briefings in bioinformatics, 6(4):357– 369, 2005.
- 107 Patrick Lewis, Myle Ott, Jingfei Du, and Veselin Stoyanov. Pretrained language models for biomedical and clinical tasks: Understanding and extending the state-of-the-art. In Proceedings of the 3rd Clinical Natural Language Processing Workshop, pages 146–157, 2020.
- 108 Huanyu Li, Zlatan Dragisic, Daniel Faria, Valentina Ivanova, Ernesto Jiménez-Ruiz, Patrick Lambrix, and Catia Pesquita. User validation in ontology alignment: functional assessment and impact. Knowl. Eng. Rev., 34:e15, 2019.
- 109 Jiao Li, Yueping Sun, Robin J. Johnson, Daniela Sciaky, Chih-Hsuan Wei, Robert Leaman, Allan Peter Davis, Carolyn J. Mattingly, Thomas C. Wiegers, and Zhiyong Lu. Biocreative V CDR task corpus: a resource for chemical disease relation extraction. Database J. Biol. Databases Curation, 2016, 2016.
- 110 Michelle M Li, Kexin Huang, and Marinka Zitnik. Graph representation learning in biomedicine and healthcare. *Nature Biomedical Engineering*, pages 1–17, 2022.
- 111Fangyu Liu, Ehsan Shareghi, Zaiqiao Meng, Marco Basaldella, and Nigel Collier. Self-alignment pretraining for biomedical entity representations. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 4228–4238, 2021.
- 112Hao Liu, Yehoshua Perl, and James Geller. Concept Placement Using BERT Trained by Transforming and Summarizing Biomedical Ontology Structure. J. of Biomedical Informatics, 112(C), 2020.
- 113Kaihong Liu, William R Hogan, and Rebecca S Crowley. Natural language processing methods and systems for biomedical ontology learning.

Journal of biomedical informatics, 44(1):163–179, 2011.

- 114Yu Liu, Jingtao Ding, Yanjie Fu, and Yong Li. Urbankg: An urban knowledge graph system. ACM Transactions on Intelligent Systems and Technology, 14(4):1–25, 2023.
- 115 Takaki Makino, Yoshihiro Ohta, Jun'ichi Tsujii, et al. Tuning support vector machines for biomedical named entity recognition. In *Proceedings of the ACL-02 workshop on Natural language processing in the biomedical domain*, pages 1–8, 2002.
- 116 Diego Marcheggiani and Ivan Titov. Discrete-State Variational Autoencoders for Joint Discovery and Factorization of Relations. Transactions of the Association for Computational Linguistics, 4:231– 244, 06 2016.
- 117 Nicolas Matentzoglu, James P. Balhoff, Susan M. Bello, Chris Bizon, Matthew H. Brush, Tiffany J. Callahan, Christopher G. Chute, William D. Duncan, Chris T. A. Evelo, Davera Gabriel, John Graybeal, Alasdair J. G. Gray, Benjamin M. Gyori, Melissa A. Haendel, Henriette Harmse, Nomi L. Harris, Ian Harrow, Harshad Hegde, Amelia L. Hoyt, Charles Tapley Hoyt, Dazhi Jiao, Ernesto Jiménez-Ruiz, Simon Jupp, Hyeongsik Kim, Sebastian Köhler, Thomas Liener, Qinqin Long, James Malone, James A. McLaughlin, Julie A. McMurry, Sierra A. T. Moxon, Monica C. Munoz-Torres, David Osumi-Sutherland, James A. Overton, Bjoern Peters, Tim E. Putman, Núria Queralt-Rosinach, Kent A. Shefchek, Harold Solbrig, Anne E. Thessen, Tania Tudorache, Nicole A. Vasilevsky, Alex H. Wagner, and Christopher J. Mungall. A simple standard for sharing ontological mappings (SSSOM). Database J. Biol. Databases Curation, 2022(2022), 2022.
- 118 Nicolas Matentzoglu, Damien Goutte-Gattat, Shawn Zheng Kai Tan, James P Balhoff, Seth Carbon, Anita R Caron, William D Duncan, Joe E Flack, Melissa Haendel, Nomi L Harris, William R Hogan, Charles Tapley Hoyt, Rebecca C Jackson, HyeongSik Kim, Huseyin Kir, Martin Larralde, Julie A McMurry, James A Overton, Bjoern Peters, Clare Pilgrim, Ray Stefancsik, Sofia MC Robb, Sabrina Toro, Nicole A Vasilevsky, Ramona Walls, Christopher J Mungall, and David Osumi-Sutherland. Ontology Development Kit: a toolkit for building, maintaining and standardizing biomedical ontologies. Database, 2022:baac087, 10 2022.
- 119 Jamie P. McCusker, Neha Keshan, Sabbir Rashid, Michael Deagen, Cate Brinson, and Deborah L. McGuinness. NanoMine: A Knowledge Graph for Nanocomposite Materials Science. In The Semantic Web – ISWC 2020: 19th International Semantic Web Conference, Athens, Greece, November 2–6, 2020, Proceedings, Part II, page 144–159, 2020.
- 120 Pierre Monnin, Miguel Couceiro, Amedeo Napoli, and Adrien Coulet. Knowledge-based matching of n-ary tuples. In Mehwish Alam, Tanya Braun, and Bruno Yun, editors, Ontologies and Concepts in Mind and Machine - 25th International Conference on Conceptual Structures, Proceedings, volume 12277 of Lecture Notes in Computer Science, pages 48–56, 2020.

- 121Pierre Monnin, Joël Legrand, Graziella Husson, Patrice Ringot, Andon Tchechmedjiev, Clément Jonquet, Amedeo Napoli, and Adrien Coulet. PGxO and PGxLOD: a reconciliation of pharmacogenomic knowledge of various provenances, enabling further comparison. *BMC Bioinform.*, 20-S(4):139:1–139:16, 2019.
- 122Pierre Monnin, Chedy Raïssi, Amedeo Napoli, and Adrien Coulet. Discovering alignment relations with graph convolutional networks: A biomedical case study. *Semantic Web*, 13(3):379–398, 2022.
- 123Deisy Morselli Gysi, Ítalo Do Valle, Marinka Zitnik, Asher Ameli, Xiao Gan, Onur Varol, Susan Dina Ghiassian, JJ Patten, Robert A Davey, Joseph Loscalzo, et al. Network medicine framework for identifying drug-repurposing opportunities for COVID-19. Proceedings of the National Academy of Sciences, 118(19):e2025581118, 2021.
- 124Boris Motik and Ljiljana Stojanovic. Ontology evolution within ontology editors. In OntoWeb-SIG3 Workshop at the 13th International Conference on Knowledge Engineering and Knowledge Management EKAW 2002; Siguenza (Spain), 30th September 2002, 2002.
- 125Lino Murali, G. Gopakumar, Daleesha M. Viswanathan, and Prema Nedungadi. Towards electronic health record-based medical knowledge graph construction, completion, and applications: A literature study. Journal of Biomedical Informatics, 143:104403, 2023.
- 126Emir Muñoz, Vít Nováček, and Pierre-Yves Vandenbussche. Facilitating prediction of adverse drug reactions by using knowledge graphs and multilabel learning models. *Briefings in Bioinformatics*, 20(1):190–202, 08 2017.
- 127Erik B Myklebust, Ernesto Jimenez-Ruiz, Jiaoyan Chen, Raoul Wolf, and Knut Erik Tollefsen. Knowledge graph embedding for ecotoxicological effect prediction. In *The Semantic Web–ISWC* 2019: 18th International Semantic Web Conference, Auckland, New Zealand, October 26–30, 2019, Proceedings, Part II 18, pages 490–506, 2019.
- 128Erik B Myklebust, Ernesto Jiménez-Ruiz, Jiaoyan Chen, Raoul Wolf, and Knut Erik Tollefsen. Prediction of adverse biological effects of chemicals using knowledge graph embeddings. *Semantic Web*, 13(3):299–338, 2022.
- 129Erik Bryhn Myklebust, Ernesto Jiménez-Ruiz, Jiaoyan Chen, Raoul Wolf, and Knut Erik Tollefsen. Understanding Adverse Biological Effect Predictions Using Knowledge Graphs. CoRR, abs/2210.15985, 2022.
- 130David Nadeau and Satoshi Sekine. A survey of named entity recognition and classification. *Lingvisticae Investigationes*, 30(1):3–26, 2007.
- 131 Yavor Nenov, Robert Piro, Boris Motik, Ian Horrocks, Zhe Wu, and Jay Banerjee. RDFox: A Highly-Scalable RDF Store. In 14th International Semantic Web Conference, volume 9367 of Lecture Notes in Computer Science, pages 3–20. Springer, 2015.
- 132Fabian Neuhaus and Janna Hastings. Ontology development is consensus creation, not (merely) representation. Applied Ontology, 17(4):495–513, 2022.

- 133David N. Nicholson and Casey S. Greene. Constructing knowledge graphs and their biomedical applications. *Computational and Structural Biotechnology Journal*, 18:1414–1428, 2020.
- 134N. Noy and D.L. McGuinness. Ontology development 101: A guide to creating your first ontology. Technical Report KSL-01-05 and SMI-2001-0880, Stanford Knowledge Systems Laboratory and Stanford Medical Informatics, 2001.
- 135 Natalya F Noy, Nigam H Shah, Patricia L Whetzel, Benjamin Dai, Michael Dorf, Nicholas Griffith, Clement Jonquet, Daniel L Rubin, Margaret-Anne Storey, Christopher G Chute, et al. Bio-Portal: ontologies and integrated data resources at the click of a mouse. *Nucleic acids research*, 37(suppl_2):W170–W173, 2009.
- 136 Vardaan Pahuja, Yu Gu, Wenhu Chen, Mehdi Bahrami, Lei Liu, Wei-Peng Chen, and Yu Su. A Systematic Investigation of KB-Text Embedding Alignment at Scale. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 1764–1774, 2021.
- 137 Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, and Xindong Wu. Unifying large language models and knowledge graphs: A roadmap. *CoRR*, abs/2306.08302, 2023.
- 138 Heiko Paulheim. Knowledge graph refinement: A survey of approaches and evaluation methods. Semant. Web, 8(3):489–508, jan 2017.
- 139 Hao Peng, Haoran Li, Yangqiu Song, Vincent Zheng, and Jianxin Li. Differentially private federated knowledge graphs embedding. In Proceedings of the 30th ACM International Conference on Information & Knowledge Management, pages 1416– 1425, 2021.
- 140 Romana Pernisch, Daniele Dell'Aglio, and Abraham Bernstein. Beware of the hierarchy—an analysis of ontology evolution and the materialisation impact for biomedical ontologies. *Journal of Web Semantics*, 70:100658, 2021.
- 141 Mina Abd Nikooie Pour, Alsayed Algergawy, Patrice Buche, Leyla Jael Castro, Jiaoyan Chen, Hang Dong, Omaima Fallatah, Daniel Faria, Irini Fundulaki, Sven Hertling, Yuan He, Ian Horrocks, Martin Huschka, Liliana Ibanescu, Ernesto Jiménez-Ruiz, Naouel Karam, Amir Laadhar, Patrick Lambrix, Huanyu Li, Ying Li, Franck Michel, Engy Nasr, Heiko Paulheim, Catia Pesquita, Tzanina Saveta, Pavel Shvaiko, Cássia Trojahn, Chantelle Verhey, Mingfang Wu, Beyza Yaman, Ondrej Zamazal, and Lu Zhou. Results of the Ontology Alignment Evaluation Initiative 2022. In Proceedings of the 17th International Workshop on Ontology Matching (OM 2022) colocated with the 21th International Semantic Web Conference (ISWC 2022), volume 3324 of CEUR Workshop Proceedings, pages 84–128, 2022.
- 142 María Poveda-Villalón, Alba Fernández-Izquierdo, Mariano Fernández-López, and Raúl García-Castro. LOT: An industrial oriented ontology engineering framework. Engineering Applications of Artificial Intelligence, 111:104755, 2022.

- 143Eric Prud'hommeaux, Steve Harris, and Andy Seaborne. SPARQL 1.1 Query Language. Technical report, W3C, 2013.
- 144Sampo Pyysalo, Filip Ginter, Juho Heimonen, Jari Björne, Jorma Boberg, Jouni Järvinen, and Tapio Salakoski. Bioinfer: a corpus for information extraction in the biomedical domain. BMC bioinformatics, 8:1–24, 2007.
- 145Zhixin Qi, Hongzhi Wang, Ziming Shen, and Donghua Yang. PreKar: A learned performance predictor for knowledge graph stores. World Wide Web, 26(1):321–341, 2023.
- 146Enayat Rajabi and Somayeh Kafaie. Knowledge graphs and explainable AI in healthcare. Inf., 13(10):459, 2022.
- 147K. E. Ravikumar, Majid Rastegar-Mojarad, Majid Rastegar-Mojarad, and Hongfang Liu. Belminer: adapting a rule-based relation extraction system to extract biological expression language statements from bio-medical literature evidence sentences. Database: The Journal of Biological Databases and Curation, 2017, 2017.
- 148KE Ravikumar, Kavishwar B Wagholikar, and Hongfang Liu. Towards pathway curation through literature mining–a case study using pharmgkb. In *Biocomputing 2014*, pages 352–363. World Scientific, 2014.
- 149 Alan Rector, Stefan Schulz, Jean Marie Rodrigues, Christopher G Chute, and Harold Solbrig. On beyond Gruber: "Ontologies" in today's biomedical information systems and the limits of OWL. *Jour*nal of Biomedical Informatics, 100:100002, 2019.
- 150Petar Ristoski and Heiko Paulheim. Semantic web in data mining and knowledge discovery: A comprehensive survey. J. Web Semant., 36:1–22, 2016.
- 151 Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle Ott, C. Lawrence Zitnick, Jerry Ma, and Rob Fergus. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. Proceedings of the National Academy of Sciences of the United States of America, 118(15), 2021.
- 152 Tim Rocktäschel, Michael Weidlich, and Ulf Leser. ChemSpot: a hybrid system for chemical named entity recognition. *Bioinformatics*, 28(12):1633– 1640, 2012.
- 153Natalia Díaz Rodríguez, Alberto Lamas, Jules Sanchez, Gianni Franchi, Ivan Donadello, Siham Tabik, David Filliat, Policarpo Cruz, Rosana Montes, and Francisco Herrera. Explainable neural-symbolic learning (X-NeSyL) methodology to fuse deep learning representations with expert knowledge graphs: The monumai cultural heritage use case. Inf. Fusion, 79:58–83, 2022.
- 154 Jerret Ross, Brian Belgodere, Vijil Chenthamarakshan, Inkit Padhi, Youssef Mroueh, and Payel Das. Large scale chemical language representations capture molecular structure and properties. *Nature Machine Intelligence*, 4(12):1256–1264, 2022.
- 155 Manuel Salvadores, Paul R. Alexander, Mark A. Musen, and Natalya Fridman Noy. BioPortal as a dataset of linked biomedical ontologies and terminologies in RDF. Semantic Web, 4(3):277–284, 2013.

- 156 Matthias Samwald, José Antonio Miñarro-Giménez, Richard D. Boyce, Robert R. Freimuth, Klaus-Peter Adlassnig, and Michel Dumontier. Pharmacogenomic knowledge representation, reasoning and genome-based clinical decision support based on OWL 2 DL ontologies. BMC Medical Informatics Decis. Mak., 15:12, 2015.
- 157 Conrad L Schoch, Stacy Ciufo, Mikhail Domrachev, Carol L Hotton, Sivakumar Kannan, Rogneda Khovanskaya, Detlef Leipe, Richard Mcveigh, Kathleen O'Neill, Barbara Robbertse, Shobha Sharma, Vladimir Soussov, John P Sullivan, Lu Sun, Seán Turner, and Ilene Karsch-Mizrachi. NCBI Taxonomy: a comprehensive update on curation, resources and tools. *Database*, 2020:baaa062, 08 2020.
- 158Lynn M Schriml, James B Munro, Mike Schor, Dustin Olley, Carrie McCracken, Victor Felix, J Allen Baron, Rebecca Jackson, Susan M Bello, Cynthia Bearer, et al. The human disease ontology 2022 update. Nucleic acids research, 50(D1):D1255–D1261, 2022.
- 159 Marta Contreiras Silva, Daniel Faria, and Catia Pesquita. Matching multiple ontologies to build a knowledge graph for personalized medicine. In The Semantic Web: 19th International Conference, ESWC 2022, Hersonissos, Crete, Greece, May 29-June 2, 2022, Proceedings, pages 461–477. Springer, 2022.
- 160 Ana Claudia Sima, Tarcisio Mendes de Farias, Erich Zbinden, Maria Anisimova, Manuel Gil, Heinz Stockinger, Kurt Stockinger, Marc Robinson-Rechavi, and Christophe Dessimoz. Enabling semantic queries across federated bioinformatics databases. *Database*, 2019:baz106, 2019.
- 161 Rita T Sousa, Sara Silva, and Catia Pesquita. Evolving knowledge graph similarity for supervised learning in complex biomedical domains. BMC bioinformatics, 21:1–19, 2020.
- 162 Rita T Sousa, Sara Silva, and Catia Pesquita. Explainable representations for relation prediction in knowledge graphs. arXiv e-prints, pages arXiv-2306, 2023.
- 163 John F Sowa et al. Semantic networks. Encyclopedia of artificial intelligence, 2:1493–1511, 1992.
- 164Lise Stork, Ilaria Tiddi, René Spijker, and Annette ten Teije. Explainable drug repurposing in context via deep reinforcement learning. In Catia Pesquita, Ernesto Jimenez-Ruiz, Jamie McCusker, Daniel Faria, Mauro Dragoni, Anastasia Dimou, Raphael Troncy, and Sven Hertling, editors, *The Semantic Web*, pages 3–20, 2023.
- 165 Kai Sun, Yuhua Liu, Zongchao Guo, and Changbo Wang. Visualization for knowledge graph based on education data. Int. J. Softw. Informatics, 10, 2016.
- 166 Mujeen Sung, Jinhyuk Lee, Sean Yi, Minji Jeon, Sungdong Kim, and Jaewoo Kang. Can language models be biomedical knowledge bases? In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 4723–4734, 2021.
- 167 Gyte Tamasauskaite and Paul Groth. Defining a knowledge graph development process through

a systematic review. ACM Trans. Softw. Eng. Methodol., 32(1), 2023.

- 168HSPO Team. Health and Social Person-centric Ontology, 9 2022.
- 169 Christos Theodoropoulos, Natasha Mulligan, Thaddeus Stappenbeck, and Joao Bettencourt-Silva. Representation learning for person or entity-centric knowledge graphs: An application in healthcare. CoRR, abs/2305.05640, 2023.
- 170Ilaria Tiddi and Stefan Schlobach. Knowledge graphs as tools for explainable machine learning: A survey. Artif. Intell., 302:103627, 2022.
- 171Santiago Timón-Reina, Mariano Rincón, and Rafael Martínez-Tomás. An overview of graph databases and their applications in the biomedical domain. *Database*, 2021:baab026, 05 2021.
- 172Igor Trajkovski, Nada Lavrac, and Jakub Tolar. SEGS: search for enriched gene sets in microarray data. J. Biomed. Informatics, 41(4):588–601, 2008.
- 173Efthymia Tsamoura, David Carral, Enrico Malizia, and Jacopo Urbani. Materializing knowledge bases via trigger graphs. *Proc. VLDB Endow.*, 14(6):943–956, 2021.
- 174Serbulent Unsal, Heval Atas, Muammer Albayrak, Kemal Turhan, Aybar C Acar, and Tunca Doğan. Learning functional properties of proteins with language models. *Nature Machine Intelligence*, 4(3):227–245, 2022.
- 175 Nicole A. Vasilevsky, Shahim Essaid, Nicolas Matentzoglu, Nomi L. Harris, Melissa A. Haendel, Peter N. Robinson, and Christopher J. Mungall. Mondo Disease Ontology: Harmonizing Disease Concepts Across the World (short paper). In Janna Hastings and Frank Loebe, editors, Proceedings of the 11th International Conference on Biomedical Ontologies (ICBO) joint with the 10th Workshop on Ontologies and Data in Life Sciences (ODLS) and part of the Bolzano Summer of Knowledge (BoSK 2020), volume 2807 of CEUR Workshop Proceedings, pages 1–2. CEUR-WS.org, 2020.
- 176Blerta Veseli, Sneha Singhania, Simon Razniewski, and Gerhard Weikum. Evaluating language models for knowledge base completion. In European Semantic Web Conference, pages 227–243, 2023.
- 177 Olga Vrousgou, Tony Burdett, Helen E. Parkinson, and Simon Jupp. Biomedical Ontology Evolution in the EMBL-EBI Ontology Lookup Service. In Themis Palpanas and Kostas Stefanidis, editors, *Proceedings of the Workshops of the EDBT/ICDT* 2016 Joint Conference, EDBT/ICDT Workshops 2016, volume 1558 of CEUR Workshop Proceedings, 2016.
- 178Hanchen Wang, Tianfan Fu, Yuanqi Du, Wenhao Gao, Kexin Huang, Ziming Liu, Payal Chandak, Shengchao Liu, Peter Van Katwyk, Andreea Deac, et al. Scientific discovery in the age of artificial intelligence. *Nature*, 620(7972):47–60, 2023.
- 179Meng Wang and Ningyu Zhang. Cross-modal knowledge discovery, inference, and challenges. In Reasoning Web. Causality, Explanations and Declarative Knowledge: 18th International Summer School 2022, Tutorial Lectures, pages 199– 209. Springer Nature Switzerland Springer, Cham, 2023.

- 180Xin Wang and Weixue Chen. Knowledge graph data management: Models, methods, and systems. In International Conference on Web Information Systems Engineering, pages 3–12. Springer, 2020.
- 181Xu Wang, Chen Yang, and Renchu Guan. A comparative study for biomedical named entity recognition. International Journal of Machine Learning and Cybernetics, 9:373–382, 2018.
- 182Xuwu Wang, Junfeng Tian, Min Gui, Zhixu Li, Rui Wang, Ming Yan, Lihan Chen, and Yanghua Xiao. WikiDiverse: A Multimodal Entity Linking Dataset with Diversified Contextual Topics and Entity Types. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 4785– 4797, 2022.
- 183M. Whirl-Carrillo, R. Huddart, L. Gong, K. Sangkuhl, C. F. Thorn, R. Whaley, and T. E. Klein. An evidence-based framework for evaluating pharmacogenomics knowledge for personalized medicine. *Clinical pharmacology and therapeutics*, 110(3):563 – 572, 2021.
- 184Xander Wilcke, Peter Bloem, and Victor De Boer. The knowledge graph as the default data model for learning on heterogeneous knowledge. *Data Sci*ence, 1(1-2):39–57, 2017.
- 185 World Health Organization. International statistical classification of diseases and related health problems. ICD-10. World Health Organization, Geneva, Switzerland, fifth edition, 2016.
- 186 Tianxing Wu, Guilin Qi, Cheng Li, and Meng Wang. A Survey of Techniques for Constructing Chinese Knowledge Graphs and Their Applications. Sustainability, 10(9):1–26, 2018.
- 187 Eryu Xia, Wen Sun, Jing Mei, Enliang Xu, Ke Wang, and Yong Qin. Mining disease-symptom relation from massive biomedical literature and its application in severe disease diagnosis. *Annual Symposium proceedings*, 2018:1118–1126, 2018.
- 188 Bo Xiong, Michael Cochez, Mojtaba Nayyeri, and Steffen Staab. Hyperbolic embedding inference for structured multi-label prediction. Advances in Neural Information Processing Systems, 35:33016– 33028, 2022.
- 189Minghao Xu, Xinyu Yuan, Santiago Miret, and Jian Tang. ProtST: Multi-Modality Learning of Protein Sequences and Biomedical Texts, 2023.
- 190Zonghai Yao, Yi Cao, Zhichao Yang, Vijeta Deshpande, and Hong Yu. Extracting biomedical factual knowledge using pretrained language model

and electronic health record context. In AMIA Annual Symposium Proceedings, volume 2022, page 1188, 2022.

- 191Ronghui You, Shuwei Yao, Hiroshi Mamitsuka, and Shanfeng Zhu. DeepGraphGO: graph neural network for large-scale, multispecies protein function prediction. *Bioinformatics*, 37(Supplement_1):i262–i271, 2021.
- 192 Ronghui You, Zihan Zhang, Yi Xiong, Fengzhu Sun, Hiroshi Mamitsuka, and Shanfeng Zhu. GO-Labeler: improving sequence-based large-scale protein function prediction by learning to rank. *Bioinformatics*, 34(14):2465–2473, 2018.
- 193 Jianbo Yuan, Zhiwei Jin, Han Guo, Hongxia Jin, Xianchao Zhang, Tristram Smith, and Jiebo Luo. Constructing biomedical domain-specific knowledge graph with minimum supervision. *Knowledge* and Information Systems, 62(1), 2020.
- 194Bohui Zhang, Albert Meroño Peñuela, and Elena Simperl. Towards explainable automatic knowledge graph construction with human-in-the-loop. In HHAI 2023: Augmenting Human Intellect, pages 274–289. IOS Press, Amsterdam, 2023.
- 195 Ningyu Zhang, Zhen Bi, Xiaozhuan Liang, Siyuan Cheng, Haosen Hong, Shumin Deng, Qiang Zhang, Jiazhang Lian, and Huajun Chen. OntoProtein: Protein Pretraining With Gene Ontology Embedding. In International Conference on Learning Representations, 2022.
- 196 Yingwen Zhao, Jun Wang, Jian Chen, Xiangliang Zhang, Maozu Guo, and Guoxian Yu. A literature review of gene function prediction by modeling gene ontology. *Frontiers in genetics*, 11:400, 2020.
- 197 Hong-Yu Zhou, Yunxiang Fu, Zhicheng Zhang, Bian Cheng, and Yizhou Yu. Protein representation learning via knowledge enhanced primary structure reasoning. In *The Eleventh International Conference on Learning Representations*, 2023.
- 198Qile Zhu, Xiaolin Li, Ana Conesa, and Cécile Pereira. GRAM-CNN: a deep learning approach with local context for named entity recognition in biomedical text. *Bioinformatics*, 34(9):1547–1554, 2018.
- 199 Marinka Zitnik, Monica Agrawal, and Jure Leskovec. Modeling polypharmacy side effects with graph convolutional networks. *Bioinformatics*, 34(13):i457–i466, 2018.
- 200Lei Zou, M Tamer Özsu, Lei Chen, Xuchuan Shen, Ruizhe Huang, and Dongyan Zhao. gStore: a graph-based SPARQL query engine. *The VLDB journal*, 23:565–590, 2014.