



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)*, 1(1), doi: 10.4230/TGDK.1.1.5

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: <https://doi.org/10.4230/TGDK.1.1.5>

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:

<http://openaccess.city.ac.uk/>

publications@city.ac.uk

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 Dong² ✉ 

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 ✉  

Université Côte d'Azur, Inria, CNRS, I3S, France

Catia Pesquita ✉

LASIGE, Faculdade de Ciências, Universidade de Lisboa, Portugal

Petr Škoda ✉ 

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 graph-structured.

The volume of data and the complexity of scientific concepts and relations referred to therein promote the application of advanced knowledge-driven 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



2012 ACM Subject Classification Information systems → Graph-based database models, Computing methodologies → Knowledge representation and reasoning, Applied computing → 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

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.

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.

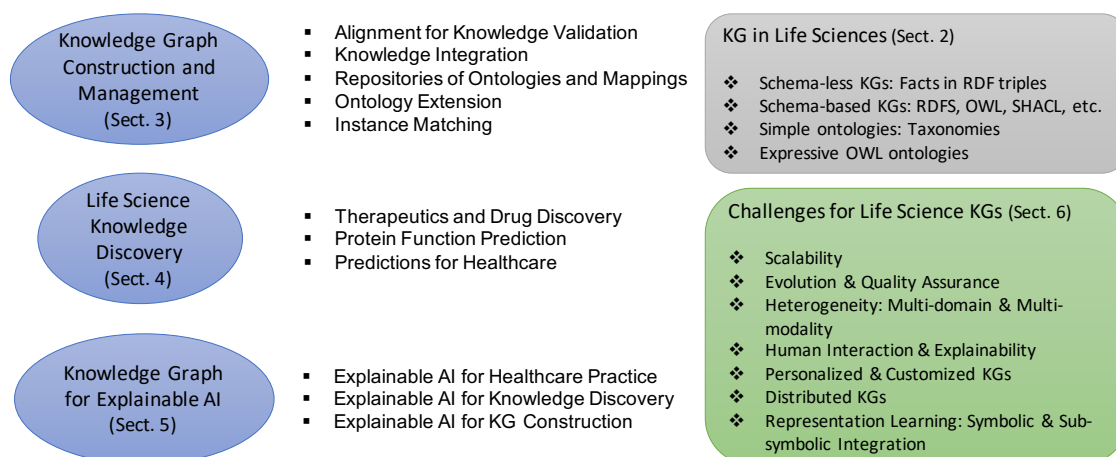
These different forms of graph-based representations can be further categorized based on the type of content represented. We therefore distinguish schema-less and schema-based Knowledge Graphs. More specifically, a typical KG contains either or both a schema part (terminologies or TBox¹) and a data part (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.

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

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

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

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

44 2 Knowledge Graphs in the Life Sciences

45 KGs represent semantically-described real-world entities, typically through ontologies (vocabu-
46 laries or schemas) [70, 63] and the data instantiating them, and thus provide descriptions of the
47 entities of interest and their interrelations, by means of links to ontology classes describing them,
48 organized in a graph [161]. KGs have been widely adopted in the life sciences, as can be seen in
49 the composition of the Linked Open Data Cloud³, where life sciences represent one of the largest
50 subdomains. A prominent example is the KG representing annotations regarding proteins by
51 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,

³ <http://cas.lod-cloud.net>

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

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

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

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

- 85 ■ Schema-less KGs composed of only relational facts in the form of RDF triples. Examples
 86 include the PharmaGKB dataset, an integrated online knowledge resource capturing how
 87 genetic variation contributes to variation in drug response [183]. Note that many semantic
 88 networks (defined in Appendix A) could be assigned to this category as their triples form a
 89 multi-relational graph.
- 90 ■ Schema-based KGs composed of relational facts and their schema (meta information) in *e.g.*,
 91 RDFS, OWL, and constraint languages such as SHACL⁶. Examples include Wikidata with its
 92 property constraints, and DBpedia with its DBpedia ontology. Whilst Wikidata and DBpedia
 93 are general-purpose KGs, they also include large-scale life science knowledge.
- 94 ■ Simple ontologies representing taxonomies. Notable examples include the tree structure of
 95 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>

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

105 **3 Knowledge Graph Construction and Management**

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

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

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

128 A recent systematic review [167] surveyed different KG development approaches to determine
129 a general development framework. The review identified six main phases that are common across
130 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>

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

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

150 3.1.1 Data source selection

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

162 3.1.2 Ontology construction

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

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

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].

180 3.1.3 Knowledge extraction

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

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

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

¹⁰<https://obofoundry.org>

220 3.1.4 Knowledge ingestion and validation

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

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

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

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

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

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

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

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

3.1.5 KG storage and inspection

KGs need to be accessible to support a variety of different tasks, beyond the mere integration of different knowledge sources, and thus KG storage management [167, 145, 180] is an active area of research. Current KG storage mechanisms are divided into relation based stores (*e.g.*, [1]) and native graph stores (*e.g.*, [200]). Relational KG stores, either based on relational databases or through NOSQL databases and / or triple stores such as Jena TDB¹¹, have reached a considerable level of maturity and have been optimized in order to avoid common problems, *e.g.*, a large number of null values in columns or optimized query performance [145]. Graph databases store nodes, edges and properties of graphs natively, and support query and graph mining tasks. Examples of state of the art implementations include Neo4J¹², GraphDB¹³, and RDFox¹⁴. The evolution of the performance of these systems has been the object of systematic studies [9], whereas [171] explicitly focuses on biomedical use cases.

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.

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

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<https://jena.apache.org/documentation/tdb/index.html>

¹²<https://neo4j.com>

¹³<https://graphdb.ontotext.com>

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

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

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

311 3.2 Examples of Life Science KG Construction

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

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

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

¹⁶<http://www.obofoundry.org/principles/fp-004-versioning.html>

¹⁷<https://biportal.bioontology.org>

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

357 3.3 What has been done: recent technical developments

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

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

¹⁸ <https://github.com/hetio/hetionet>

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

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.

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

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

²⁰ <https://pubchem.ncbi.nlm.nih.gov/>

²¹ <https://www.ebi.ac.uk/chebi/>

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

438 3.4 What are the challenges?

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

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

479 **4 Life Science Knowledge Discovery**

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

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

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

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

511 **4.1 What has been done: use cases and their recent developments**

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

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

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

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

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

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

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

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.

- 601 ■ **How to incorporate the semantics from a KG in machine learning?** Many life
602 science knowledge discovery tasks are modeled as a machine learning classification problem,
603 whose input and output labels have additional valuable information in one or multiple ex-
604 ternal KGs. The challenge lies in extracting this information, optionally encoding it into
605 vector representations, and injecting that knowledge into machine learning and pre-trained
606 language models. Doing this effectively remains an important open challenge especially for
607 protein-related pre-trained language models [195, 189, 197]. Besides improving the accuracy
608 in knowledge discovery, injecting semantics from KGs can also contribute to making the model
609 more explainable (see Section 5), but to this end, much research is still required.
- 610 ■ **How to deal with the long-tail phenomenon in machine learning with KGs?** In
611 machine learning classification for real-world life science knowledge discovery, the candidate
612 labels often exhibit a long-tailed distribution, *i.e.*, a small ratio of them are common with a
613 large number of training samples available, while most of them are infrequent or even have
614 never appeared before. For example, imbalance in training data may occur for rare diseases
615 or adverse drug effects that affect only a small portion of the population [2, 74, 38]. KGs
616 sometimes have encoded the relationships of the labels, and could be used to help train the
617 model for predicting those long-tailed labels or enable the inference of such labels.

- 618 ■ **How to create an efficient multi-modal representation of knowledge to enable dis-**
619 **covery?** Most current state-of-the-art methods build learned graph representations based on
620 isolated modalities. Multimodal KGs can explicitly capture labelled nodes and edges, each
621 with well-defined meanings, across heterogeneous node types, relations and modalities (such
622 as text, images, protein sequences, molecules fingerprints, diseases and more) [20, 103]. Incorporating
623 KGs with multiple modalities for representation learning requires computationally
624 scalable methods to compute the initial embeddings for each modality, as a preliminary step
625 to learn computable representations of large knowledge. Furthermore, robust learning tech-
626 niques are needed for generalizing the learned representations to nodes with unseen or missing
627 modalities, thereby enabling the discovery of new knowledge. An example would be inferring
628 properties of proteins for which only the sequence is known.
- 629 ■ **How to efficiently utilize and fuse heterogeneous datasets, such as human-curated**
630 **domain knowledge bases, scientific literature and person-centered health records,**
631 **for knowledge discovery?** State of the art shows that representations can be enhanced by
632 incorporating richer information available across different sources [74, 103, 159]. Bringing in
633 more data during training is needed to learn representations that can be applied to a broader
634 range of downstream prediction tasks. However, learning from large and diverse KGs requires
635 addressing challenges such as alignment, noise handling, balancing rich expressive knowledge
636 with scalability and dealing with knowledge inconsistency. Moreover, more robust learning
637 methods are needed for generalizing the learned representation to multiple downstream tasks
638 (*e.g.*, knowledge-aware transfer, zero-shot and few-shot learning [23]). An important aspect
639 in this regard is addressing the disparity between all of the knowledge accessible during pre-
640 training and the knowledge accessible or relevant for the downstream fine-tuning [74, 103].

641 **5 Knowledge Graphs for Explainable AI**

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

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

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

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

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

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

686 5.1 What has been done: use cases and recent developments

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

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

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

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

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

730 5.2 What are the challenges?

- 731 ■ **How to integrate KGs for better XAI, especially with recent deep learning and**
732 **language model based methods?** KG may provide better data provenance for the model
733 output. This can ensure explainability for communicating the model to domain experts in
734 data science applications [8]. In terms of recent generative LLMs, life science KGs, with careful
735 curation based on scientific publications, may help to provide provenance data to the answers
736 generated by LLMs. Studies need to understand to what extent, and how, LLMs can be applied
737 to induce knowledge (*e.g.*, by probing LLMs with biomedical ontologies [67]), which then may
738 provide a foundation to create better approaches to integrate KGs with LLMs. Another area
739 is neuro-symbolic methods which may provide models that are inherently more interpretable
740 (see further discussions in Section 6.1). Also, regarding language models (especially LLMs),
741 they are capable of generating fluent texts, which can potentially serve as textual explanation
742 generators from symbolic knowledge for XAI. Meanwhile, a key issue is the hallucination of
743 LLMs, and KGs may support better prompting, fine-tuning and interpretable inference of
744 LLMs for higher decisiveness and trustfulness [137].
- 745 ■ **How to evaluate XAI methods that involve KG?** How to measure the quality of ex-
746 planations, to ensure they are corresponding to users? The majority (around 70%) of XAI
747 studies for KG construction do not evaluate the quality of the explanations or only informally
748 visualize or comment on a limited number of cases to show the intuitive outcome [194]. Also,
749 an XAI method needs to consider the target audience, as the explainability is to be finally
750 received by a group of humans [8]. For instance, only a small number of current approaches
751 to XAI for KG construction involve a user study, human evaluation or task-specific met-
752 rics [194]. Evaluating the quality of explanations requires some expert evaluation performed
753 as ex-post evaluation, and well-defined metrics are needed for this task. An example is in [59]
754 to use a combination of users' scores for each predicted explanation in a KG link prediction
755 task, where there are multiple possible explanations. More expert validated and automated

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

757 **6 Discussion and Conclusion**

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

762 **6.1 Overall challenges and trends**

763 Meanwhile, more scalable and efficient knowledge retrieval, query and reasoning systems, includ-
764 ing life science KGs and mapping repositories, are still worthy of investigation and development.

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

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

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

800 [92] summarizes XAI in bioinformatics with a chapter on knowledge-based explanations, whereas
801 Karim [91, Chapter 8] provides a neuro-symbolic framework for KG construction and utilisation
802 for medical experts' decision making in the cancer domain. The approach presented in [153] is
803 another recent example of neuro-symbolic integration for image classification with KG-based XAI
804 in the cultural heritage domain.

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

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

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

844 **6.2 Conclusion**

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

857 **Appendix A: Terms in Knowledge Graphs and Life Sciences**

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

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

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

871 **Semantic Networks:** a graph structure for representing knowledge in patterns of intercon-
 872 nected nodes and arcs [163]. We use the term to denote a graph of concepts and relations without
 873 formal semantics.

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

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

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

887 **ChEBI:** Chemical Entities of Biological Interest (ChEBI) is a database and ontology contain-
 888 ing information about chemical entities of biological interest. [64]

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

tems of human cognition are related to the representation and manipulation of symbols; sub-symbolic 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 using neural network approaches to predict words from prior their contexts in a sequence. Pre-training is the process of learning some sort of representation (usually neural embedding based) of meaning for words or sentences by processing very large amounts of text (or other data in a sequence form, *e.g.*, proteins and KG facts). This results in pre-trained language models. The dominating architecture for neural language modeling is Transformer-based models, including BERT, its domain specific versions, and later large variants, like the GPT series. The pre-trained language models of very large sizes are recently coined Large Language Models (LLMs). [88]

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]

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

- 1 Daniel 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.
- 2 Emily 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.
- 3 Robert Arp, Barry Smith, and Andrew D. Spear. *Building Ontologies With Basic Formal Ontology*. The MIT Press, 08 2015.
- 4 M 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.
- 6 Franz Baader, Ian Horrocks, Carsten Lutz, and Uli Sattler. *An Introduction to Description Logic*. Cambridge University Press, Cambridge, 2017.
- 7 Nguyen Bach and Sameer Badaskar. A review of relation extraction. *Literature review for Language and Statistics II*, 2:1–15, 2007.
- 8 Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bénézet, Siham Tabik, Alberto Barbado, Salvador García, 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 Hoeffler. Demystifying graph databases: Analysis and taxonomy of data organization, system designs, and graph queries. *ACM Comput. Surv.*, 2023.
- 10 Olivier Bodenreider. The Unified Medical Language System (UMLS): integrating biomedical terminology. *Nucleic Acids Res.*, 32(Database-Issue):267–270, 2004.
- 11 Nadav 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.
- 13 Emmanuel Bresso, Pierre Monnin, Cédric Bousquet, François-Élie Calvier, Ndeye Coumba Ndiaye, Nadine Petitpain, Malika Smail-Tabbone, and Adrien Coulet. Investigating ADR mechanisms with explainable AI: a feasibility study with knowledge graph mining. *BMC Medical Informatics Decision Mak.*, 21(1):171, 2021.
- 14 Tom 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.
- 15 Rich Caruana, Yin Lou, Johannes Gehrke, Paul Koch, Marc Sturm, and Noemie Elhadad. Intelligent 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 Matentzoglou, HyeonSik 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. Low-dimensional hyperbolic knowledge graph embeddings. *arXiv preprint arXiv:2005.00545*, 2020.
- 20 Payal Chandak, Kexin Huang, and Marinka Zitnik. Building a knowledge graph to enable precision medicine. *bioRxiv*, 2022.
- 21 Shruthi 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 Jiaoyan 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.
- 28 E. 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 Alexander, 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 Asanithong, Diana Luna Buitrago, Meltem N Erdol, Matthew C Gage, Mohamed Ali Kadhum, Kan Yan Chloe Li, Miao Long, Aleksandra Michalak, Angeline

- Pesala, Armalya Pritazahra, Shirin C C Saverimuttu, Renzhi Su, Kate E Thurlow, Ruth C Lovering, Colin Logie, Snezhana Oliferenko, Judith Blake, Karen Christie, Lori Corbani, Mary E Dolan, Harold J Drabkin, David P Hill, Li Ni, Dmitry Sitnikov, Cynthia Smith, Alayne Cuzick, 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 Streicher, 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 Ish-tiaq, 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, Leyla Ruzicka, and Monte Westerfield. The Gene Ontology knowledge-base in 2023. *Genetics*, 224(1):iyad031, 03 2023.
- 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.
- 32 Shumin 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 billion-scale 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.
- 34 Ivan Donadello, Mauro Dragoni, and Claudio Echer. Persuasive explanation of reasoning inferences on dietary data. In *PROFILES/SEMEX@ISWC*, 2019.
- 35 Hang 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.
- 36 Hang Dong, Jiaoyan Chen, Yuan He, Yinan Liu, and Ian Horrocks. Reveal the unknown: Out-of-knowledge-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.
- 38 Hang 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.
- 39 Kevin Donnelly et al. SNOMED-CT: The advanced terminology and coding system for ehealth. In *Medical and Care Compunetics 3*, volume 121 of *Studies in health technology and informatics*, pages 279–290. IOS Press, 2006.
- 40 Kevin Donnelly et al. SNOMED-CT: The advanced terminology and coding system for ehealth. In *Medical and Care Compunetics 3*, volume 121 of *Studies in health technology and informatics*, pages 279–290. IOS Press, 2006.
- 41 Zlatan 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.
- 42 Juan M Durán. Dissecting scientific explanation in ai (sxai): A case for medicine and healthcare. *Artificial Intelligence*, 297:103498, 2021.
- 43 Artur 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.
- 45 Lisa Ehrlinger and Wolfram WöB. Towards a definition of knowledge graphs. In *International Conference on Semantic Systems*, 2016.
- 46 Daniel 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 Vrandečić, Paul Groth, Natasha F. Noy, Krzysztof Janowicz, and Carole A. Goble, editors, *The Semantic Web - ISWC 2014 - 13th International Semantic Web*

- 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 Gaudet, 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.
- 53 David 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.
- 54 Aryo 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](https://arxiv.org/abs/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 Ruggeri, Franco Turini, Fosca Giannotti, and Dino Pedreschi. A survey of methods for explaining black box models. *ACM Comput. Surv.*, 51(5), aug 2018.
- 58 Ricardo 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.
- 59 Nicholas 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.
- 60 I. 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.
- 64 Janna 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.
- 68 Nicolas 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.
- 69 Daniel 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.
- 71 Aidan 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, Sabir 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.
- 74 Kexin 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.
- 75 Kexin 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 Hugué 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.
- 77 L Hunter, Z Lu, and J Firby. Wab jr, hl johnson, pv ogren, and kb cohen, “opendmap: An open source, ontology-driven concept analysis engine, with applications to capturing knowledge regarding protein transport, protein interactions and cell-type-specific 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.
- 79 Nicolas 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.
- 81 Shaoxiong 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.
- 82 Ernesto 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.
- 83 Ernesto Jiménez-Ruiz, Bernardo Cuenca Grau, Ian Horrocks, and Rafael Berlanga Llavori. Logic-based 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>.
- 84 Ernesto 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.
- 85 Ernesto 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 Vrousitou, 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.
- 87 Simon 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 International 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.
- 92 Md. 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.
- 93 Md. 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.
- 98 Elisa F. Kendall and Deborah L. McGuinness. *Ontology Engineering*. Synthesis Lectures on the Semantic Web: Theory and Technology. Springer, Cham, Switzerland, 2019.
- 99 Sebastian Köhler, Michael Gargano, Nicolas Matentzoglou, 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.
- 102 Maxat Kulmanov and Robert Hoehndorf. DeepGOZero: improving protein function prediction from sequence and zero-shot learning based on ontology axioms. *Bioinformatics*, 38(Supplement_1):i238–i245, 2022.
- 103 Hoang 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.
- 104 Robert 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. Wieggers, 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.
- 111 Fangyu Liu, Ehsan Shareghi, Zaiqiao Meng, Marco Basaldella, and Nigel Collier. Self-alignment pre-training 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.
- 112 Hao 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.
- 113 Kaihong 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.
- 114 Yu 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 Matentzoglou, 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, Qin-qin 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 Matentzoglou, 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 Laralde, 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.
- 121 Pierre 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.
- 122 Pierre 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.
- 123 Deisy 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.
- 124 Boris 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.
- 125 Lino 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.
- 126 Emir Muñoz, Vít Nováček, and Pierre-Yves Vandebussche. Facilitating prediction of adverse drug reactions by using knowledge graphs and multi-label learning models. *Briefings in Bioinformatics*, 20(1):190–202, 08 2017.
- 127 Erik 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.
- 128 Erik 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.
- 129 Erik 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.
- 130 David Nadeau and Satoshi Sekine. A survey of named entity recognition and classification. *Linguisticae 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.
- 132 Fabian Neuhaus and Janna Hastings. Ontology development is consensus creation, not (merely) representation. *Applied Ontology*, 17(4):495–513, 2022.

- 133 David N. Nicholson and Casey S. Greene. Constructing knowledge graphs and their biomedical applications. *Computational and Structural Biotechnology Journal*, 18:1414–1428, 2020.
- 134 N. 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. BioPortal: 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, Wenhui 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’Aglia, 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, Omaira 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) co-located 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.
- 143 Eric Prud’hommeaux, Steve Harris, and Andy Seaborne. SPARQL 1.1 Query Language. Technical report, W3C, 2013.
- 144 Sampo 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.
- 145 Zhixin 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.
- 146 Enayat Rajabi and Somayeh Kafaie. Knowledge graphs and explainable AI in healthcare. *Inf.*, 13(10):459, 2022.
- 147 K. 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.
- 148 KE Ravikumar, Kavishwar B Wagholikar, and Hongfang Liu. Towards pathway curation through literature mining—a case study using pharmsgkb. 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. *Journal of Biomedical Informatics*, 100:100002, 2019.
- 150 Petar 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, Sidharth 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.
- 153 Natalia 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 Chenthamarashan, 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 Ciuffo, 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.
- 158 Lynn 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.
- 164 Lise 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 Tamasauskaitė and Paul Groth. Defining a knowledge graph development process through a systematic review. *ACM Trans. Softw. Eng. Methodol.*, 32(1), 2023.
- 168 HSPO 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.
- 170 Ilaria Tiddi and Stefan Schlobach. Knowledge graphs as tools for explainable machine learning: A survey. *Artif. Intell.*, 302:103627, 2022.
- 171 Santiago 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.
- 172 Igor Trajkovski, Nada Lavrac, and Jakub Tolar. SEGs: search for enriched gene sets in microarray data. *J. Biomed. Informatics*, 41(4):588–601, 2008.
- 173 Efthymia Tsamoura, David Carral, Enrico Malizia, and Jacopo Urbani. Materializing knowledge bases via trigger graphs. *Proc. VLDB Endow.*, 14(6):943–956, 2021.
- 174 Serbulent 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 Matentzoglou, 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.
- 176 Blerta 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 Vrousou, 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.
- 178 Hanchen 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.
- 179 Meng 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.

- 180 Xin 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.
- 181 Xu 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.
- 182 Xuwu 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.
- 183 M. 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.
- 184 Xander Wilcke, Peter Bloem, and Victor De Boer. The knowledge graph as the default data model for learning on heterogeneous knowledge. *Data Science*, 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.
- 189 Minghao Xu, Xinyu Yuan, Santiago Miret, and Jian Tang. ProtST: Multi-Modality Learning of Protein Sequences and Biomedical Texts, 2023.
- 190 Zonghai 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.
- 191 Ronghui 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.
- 194 Bohui Zhang, Albert Meroño Peñuela, and Elena Simperl. Towards explainable automatic knowledge graph construction with human-in-the-loop. In *HAI 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.
- 198 Qile 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.
- 200 Lei 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.