

Towards Explainable Automated Knowledge Engineering with Human-in-the-loop

Bohui Zhang^{a,*}, Albert Meroño-Peñuela^a and Elena Simperl^a

^a*Department of Informatics, King's College London, London, United Kingdom*

E-mails: bohui.zhang@kcl.ac.uk, albert.merono@kcl.ac.uk, elena.simperl@kcl.ac.uk

Abstract. Knowledge graphs are important in human-centered AI as they provide large labeled machine learning datasets, enhance retrieval-augmented generation, and generate explanations. However, knowledge graph construction has evolved into a complex, semi-automatic process that increasingly relies on black-box deep learning models and heterogeneous data sources to scale. The knowledge graph lifecycle is not transparent, accountability is limited, and there are no accounts of, or indeed methods to determine, how fair a knowledge graph is in downstream applications. Knowledge graphs are thus at odds with AI regulation, for instance, the EU's AI Act, and with ongoing efforts elsewhere in AI to audit and debias data and algorithms. This paper reports on work towards designing explainable (XAI) knowledge-graph construction pipelines with humans in-the-loop and discusses research topics in this area. Our work is based on a systematic literature review, in which we study tasks in knowledge graph construction that are often automated, as well as common methods to explain how they work and their outcomes, and an interview study with 13 people from the knowledge engineering community. To analyze the related literature, we introduce use cases, their related goals for XAI methods in knowledge graph construction, and the gaps in each use case. To gain an understanding of the role of explainable models in practical scenarios, and reveal the requirements for improving the current XAI methods, we designed interview questions covering broad transparency and explainability topics, along with example discussion sessions using examples from the literature review. From practical knowledge engineering experience, we collect requirements for designing XAI methods, propose design blueprints, and outline directions for future research: (i) tasks in knowledge graph construction where manual input remains essential and where AI assistance could be beneficial; (ii) integrating XAI methods into established knowledge engineering practices to improve stakeholder experience; (iii) the need to evaluate how effective explanations genuinely are making human-machine collaboration in knowledge graph construction more trustworthy; (iv) adapting explanations for multiple use cases; and (v) verifying and applying the XAI design blueprint in practical settings.

Keywords: knowledge graph, knowledge graph construction, knowledge engineering, transparency, explainability, explainable AI, trustworthy AI

1. Introduction

To reach its potential, AI needs data and context. Without the right (amounts of) data, machine learning (ML) cannot identify patterns or make predictions. Without a deeper understanding of context, AI applications cannot engage people in a meaningful way. Knowledge graphs (KGs) [1], a term coined by Google in 2012 to refer to its general-purpose knowledge base, are critical to both: they reduce the need for large labeled ML datasets [2], enhance pre-trained language models (PLMs) [3, 4], and generate explanations [5]. KGs are routinely used alongside ML in many applications, including search, question answering, recommendation [1] and, in industry contexts, enterprise data management, digital twins, supply chain management, procurement, and regulatory compliance [6]. Moreover,

*Corresponding author. E-mail: bohui.zhang@kcl.ac.uk.

with the rise of large language models (LLMs) such as GPT [7, 8] and Llama series [9, 10], KGs and LLMs have influenced each other in both ways: LLMs for KGs (using LLMs for KG construction and maintenance) and KGs for LLMs (using KGs to train, prompt, augment, and evaluate LLMs) [11].

As AI applications produce and consume more data, engineering KGs has evolved into a complex, semi-automatic process that increasingly relies on opaque deep-learning models and vast collections of heterogeneous data sources to scale to graphs with millions of entities and billions of statements [12, 13]. The KG lifecycle is not transparent [14], accountability is limited, and accounts of how biased a KG is [15] or how fair the downstream applications that use it are patchy [16]. KGs are thus at odds with AI regulation, for instance, the UK's AI regulation¹, the EU's AI Act², and the US's AI Risk Management Framework (AI RMF 1.0)³ with ongoing efforts elsewhere in AI to systematically audit and debias data and algorithms and to enhance AI trustworthiness [17–21]. Most regulators take a risk-based approach to the use of AI, prescribing, among other things, transparency and accountability obligations for different classes of AI applications. Organizations using KGs, either directly as data infrastructure, or as graph embeddings in ML models, need to document and attest that their KGs are compliant with the law. Furthermore, when a KG is part of an AI application that counts as high-risk, that application will have to undergo conformity assessments both at design and at run time. KGs themselves are meant to make ML models explainable [5] and hence facilitate such compliance tasks, but that would imply that the KG lifecycle abides by the same rules.

We argue that this is not yet the case. As referred to in our previous work [22, 23], questions regarding the user-centric aspects of knowledge engineering are not yet fully answered, such as users' tasks and goals, the way that they interact with KGs, KG construction tools, and KG-related applications. Up-to-date comparative surveys regarding the scale, complexity, and degree of automation of knowledge graph construction systems nowadays are needed. User-centric design and empirical methods should be established for transparent knowledge graph construction to ensure that human-centric challenges are not overlooked.

With this paper, we would like to advance the field of **explainable knowledge engineering** to allow KG stakeholders to rely appropriately on AI algorithms and use KGs with confidence [24]. We need to first gain a better understanding of emerging KG construction practices in the era of ML-as-a-service and develop human-in-the-loop approaches to ensure transparency and accountability throughout the KG lifecycle. This applies both to proprietary KGs used within organizations [6] and publicly available KGs like Wikidata [25], DBpedia [26], YAGO [27], and ConceptNet [28], which are extensively used by researchers and practitioners. As AI laws and regulations enter into force, the trustworthy credentials of such KGs will have to be systematically assessed and documented.

Our paper follows recent work that explores emergent neuro-symbolic AI architectures from a system-design perspective. Van Bekkum et al. [29] propose a taxonomy of hybrid (i.e., learning and reasoning) systems and discuss common architecture patterns and use cases. Building on their insights, Breit et al. [30] carried out a comprehensive literature review to add details to those patterns in terms of inputs, outputs, processing units, types of ML models and their training, types of knowledge representation and reasoning, but also transparency and auditability. One of their main findings is that most system designers do not consider these latter aspects at all, or, when they do, they do not evaluate them sufficiently. A third paper by Tamašauskaitė and Groth [12] draw from a survey of system papers to define a canonical KG construction process. Our work continues where they left off: starting from their KG construction process, we follow one of their main recommendations to map tools and techniques for each step to provide additional guidance to researchers and developers.

Thus, we put forth the following research questions:

- **RQ 1:** What is the state-of-the-art/status of explainable automated knowledge graph construction?
- **RQ 2:** How do knowledge engineers and knowledge graph researchers understand their models and tools and explain their output to stakeholders?
- **RQ 3:** Do the existing explainable models and tools meet the requirements of knowledge engineers in practical use cases?

¹<https://www.gov.uk/government/publications/ai-regulation-a-pro-innovation-approach>

²<https://artificialintelligenceact.eu>

³<https://doi.org/10.6028/NIST.AI.100-1>

1 – **RQ 4:** What are the requirements of knowledge engineers and knowledge graph researchers for explainable
2 approaches?
3

4 We analyze the KG lifecycle to identify tasks that are commonly automated with AI and those that still require
5 human input and oversight and could potentially benefit from AI assistance. In parallel, we survey the state-of-the-art
6 in explainable AI (XAI) to inform the design of XAI approaches that are genuinely useful for KG stakeholders such
7 as knowledge engineers, subject domain experts, and users. We maintain a public repository to enhance research
8 convenience⁴. Then we conduct an interview study involving 13 knowledge engineers and researchers from the
9 knowledge engineering community. The interviews further explore topics such as their degree of understanding of
10 methodologies and tools, their degree of automation, their transparency and explainability requirements, and various
11 usage scenarios. Our main findings are:
12

- 13 1. There are tasks in KG construction, for instance, knowledge acquisition, where automation⁵ is routinely used
14 with promising results. At the same time, there are opportunities to use AI to assist other tasks, including
15 ontology reuse, ontology evolution, ontology evaluation, and documentation, where (the latest) AI capabilities
16 have remained under-explored.
17
- 18 2. While tasks around knowledge acquisition, taxonomy building, and data ingestion are often automated, human
19 oversight is still needed to improve performance, establish trust, or comply with the law. In our review, we
20 found little evidence of the integration of AI capabilities besides basic automation, no matter their level of
21 interpretability, into standard knowledge-engineering tools and practices. Furthermore, our understanding of
22 human-in-the-loop KG construction remains limited, with implications for user experience.
23
- 24 3. Comprehensive evaluations of XAI methods are lacking, with most studies focusing on simple ML models in
25 lab settings, with mixed results [31–33]. The KG community, just like elsewhere in AI, needs to gain a better
26 understanding of how people react and use explanations to build trust and boost technology adoption.
27
- 28 4. Knowledge engineers have varying levels of understanding regarding the tools and models they use, with many
29 expressing concerns over the opaqueness of black-box models. Data provenance and lineage tracking are rec-
30 ognized as critical, yet there are still gaps in the comprehensiveness and standardization of these practices.
31 Evaluation heavily relies on human effort, highlighting the need for more robust and scalable methods. Addi-
32 tionally, effective communication of tool functionality and results to diverse stakeholders remains a significant
33 challenge, requiring tailored approaches to bridge knowledge gaps and align expectations.
34
- 35 5. Several use cases for XAI models, such as understanding performance and contributing factors, model debug-
36 ging, enhancing human-machine interactions, and uncovering new and previously unnoticed insights. How-
37 ever, participants found that current XAI solutions often fall short of practical requirements. Issues include
38 explanations being insufficiently informative, overly complex, and lacking stability and coverage. Addition-
39 ally, participants emphasized the need for explanations to be both clear and confidence-indicating, with a
40 strong preference for natural language representations.

41 Based on these findings, we propose several directions for future research, drawing on theory and insights from
42 AI, human-AI interaction [34], interactive ML [35], and social computing [36, 37]. These include: (i) AI assistants
43 for overlooked tasks in the KG lifecycle; (ii) end-to-end tools supporting automated KG construction with human-
44 in-the-loop with built-in advanced, explainable AI capabilities; (iii) holistic evaluation frameworks that assess the
45 extent to which explanations genuinely help humans engineer better KGs; (iv) explanations adaptable for multiple
46 use cases; and (v) applying and verifying the XAI design blueprint in practical settings.
47
48

49 ⁴<https://github.com/bohuizhang/XKGC>

50 ⁵In this paper we use AI assistance and automation interchangeably. While we acknowledge that not all automation in KG construction is AI,
51 we argue that the use of AI brings about specific challenges with respect to transparency, accountability etc.

2. Background

2.1. Transparency and Explainability of ML Methods

Transparency as an AI design principle stands for the need to clearly document and explain how an AI system makes decisions, how the data is collected, used, and governed, and how the system is evaluated and audited [38–40]. Achieving transparency in machine learning (ML) models can be accomplished through explainability. Although some ML models, like decision trees, are naturally interpretable, larger models, such as language models, are too complex to comprehend in the same way. To address this issue, researchers and practitioners have proposed many XAI frameworks, guidance, standards [41], techniques [42, 43], and evaluation metrics [44] for various models within the context of trustworthy AI. Reviews and surveys of emerging XAI methods and ever-changing intents and requirements from end-users have also become increasingly common. Typically, XAI surveys [41, 45–48] focus on aspects like problem formulation, taxonomies and classification, evaluation metrics, challenges, and future directions. For works that are more related to ours, Danilevsky et al. [49] conducted a survey on the state-of-the-art XAI models in natural language processing, which includes tasks that overlap with our work, such as named entity recognition and relation extraction. In the area of XAI and KGs, researchers have suggested using KGs to provide explanations. Tiddi and Schlobach’s systematic literature review [5] focused on the integration of KGs into explainable machine learning, where KGs are used as domain knowledge for explanations. In addition to the technical perspective, Miller’s review [50] provided a thorough examination of explainable AI through a sociotechnical lens, drawing from a variety of fields such as philosophy, cognitive science, and social psychology. Although previous studies have focused on some knowledge graph construction tasks and applications, a thorough review of the transparency and explainability of knowledge graph construction is still missing.

2.2. User Studies on Explainable AI

A deep understanding of the end-user requirements is essential in order to design trustworthy explanations, as explainability is a human-centric property [51]. Preece et al. [52] give an analysis of stakeholders in XAI by examining the concerns of various stakeholders communities and digging into their different intents and requirements. Ras et al. distinguished different users of deep learning models into two groups and discussed their concerns: the expert users, who are engineers and developers building and maintaining the systems, and lay users, who are the end users and stakeholders [53]. Liao et al. [54] conducted interviews with UX and design practitioners working on various AI products through question-driven explanations. It is noteworthy that there is a lack of user studies on XAI involving knowledge engineers and knowledge graph stakeholders as end-users. Therefore, there is no consensus among design disciplines for XAI in relevant domains. Similar to our intents, Dhanorkar et al. [55] conducted an interview study on XAI towards AI researchers and stakeholders in industrial AI projects focusing on the AI lifecycle. Rong et al. [51] surveyed user studies through characteristics including trust, fairness, understanding, usability, and human-AI collaboration performance, and provided guidelines for both XAI researchers and practitioners on designing and conducting user studies. Similar to our interview study, Kim et al. [56] conducted an interactive feedback session in their interview study with the objective of understanding how explainability can support human-AI interaction. They mock up explanations that could be potentially used for AI application outputs in the field of computer vision to assess the participant’s perception of existing XAI approaches and how participants use explanations during their collaboration with the AI. Automated and transferable evaluation, benchmarking, and comparison of XAI approaches pose open challenges, as explainability is often seen as a subjective property, necessitating auditing from multiple aspects [57]. On the other hand, human-centered XAI evaluations that take an HCI perspective remain critical in XAI evaluation, where rigorous evaluation procedures need to be established [58].

2.3. Human-Centric Knowledge Engineering

Knowledge engineering, the branch of AI concerned with building and managing knowledge-based systems [59, 60], has changed dramatically with the latest innovations in machine learning, natural language processing,

1 and computer vision. The process of constructing a knowledge graph can take on various forms, but it usually involves acquiring knowledge, processing it, and deploying the knowledge graph [1, 12, 61]. And yet, as the most recent advances in natural language processing (especially LLMs) and generative AI demonstrate, the question of how to capture and encode domain knowledge into a computational representation remains as challenging as ever [62]. The technologies and end-user tools to support core knowledge-engineering tasks such as knowledge acquisition have advanced significantly to meet the scale requirements of modern KGs and to leverage the generative ability of sequence-to-sequence frameworks [63, 64]. AI copilots, which leverage LLMs, have also become involved in the KG lifecycle through conversational interactions [65], assisting knowledge engineers and users in a wide range of tasks. At the same time, the most effective approaches to knowledge representation still require human oversight at various levels [66, 67], but increasingly human input is in the form of enhancing or validating algorithmic suggestions [12]. The tasks of knowledge engineering require human-in-the-loop to a different extent and are considered human-centric [23, 68, 69]. These developments have resulted in improved methods and tools to support the knowledge engineering process, with a growing group of participants and stakeholders, including knowledge engineers and domain experts [66]. Witschel et al. identified human-in-the-loop patterns in hybrid learning and knowledge engineering activities, encapsulating them in two boxologies, where human agents function either as feedback-providers or feedback-consumers [69]. Back to 2002, after Holsapple and Joshi introduced the first collaborative approach to ontology design [70], various collaborative ontology engineering methodologies have been proposed, including tasks like ontology design and construction [71–74], ontology evolution [71, 74–76], and ontology evaluation [77, 78]. The tasks of ontology engineering continue to rely heavily on manual labor, and many of the reviewed works are outdated and pre-date the era of deep learning. There are evident challenges in improving the methodologies used in this process and adapting them to meet the requirements of automation, scalability, and transparency.

2.4. The KG Lifecycle

Building on the process from [12], Figure 1 shows that the KG lifecycle today consists of four stages with a mix of automated and manual capabilities and contributions from several stakeholder groups: knowledge engineering and machine learning specialists, subject domain experts, online volunteers, and crowdsourcing services, as well as developers of applications using KGs.

As the figure suggests, **KGs are interacting with AI capabilities in complex ways**. Human-in-the-loop tasks in KG lifecycle increasingly use ML models with varying levels of interpretability. On the left side of the figure, at stage A, which is an entry point and essential step of the KG lifecycle, knowledge engineers and KG stakeholders (e.g., domain experts) will first determine the scope of work and the success criteria [79]. After that, at the second stage, knowledge graph construction, knowledge engineers and other specialists (potentially) reuse standard ontologies and build knowledge graphs from scratch through data lifting and knowledge extraction. Multiple data sources, structured and unstructured, are lifted into KGs using ML for named entity recognition [80], relation extraction [81], entity reconciliation [82], and many others. The ontology organizing the KG can be provided upfront or derived from the data itself, depending on whether there is a clear domain or available structured data with predefined types of entities and relations [12]. In this context, [14] discusses the need for more transparency with respect to data provenance and currency; both can affect whether application developers and end-users will be able to use the KG with confidence as a source of reliable, complete, unbiased, and up-to-date information. KGs can also be created on a larger scale through human collaboration, utilizing crowdsourcing platforms, collaborative-editing platforms, etc [1]. Crowd workers and volunteer editors have important roles in the KG lifecycle, especially in knowledge graph creation and updates, where annotation tasks such as quizzes and voting are often designed for leveraging their background knowledge [83–85]. While KGs constructed using these approaches may exhibit quality issues such as errors [86, 87], disagreement [88], bias [1], etc., crowdsourcing for supervised ML has similar transparency challenges as the algorithms it complements. This is because the digital services commonly used for this purpose, e.g., Prolific and Mechanical Turk, are black-box, proprietary platforms with limited means to replicate or reproduce results [89]. Educating crowd workers in the process of performing crowdsourcing tasks is also a nontrivial task [84]. Interleaving explanations during this process could aid in educating crowd workers, enhancing their comprehension of the task, and ultimately improving output quality.

The result of knowledge acquisition is shown in the middle of the figure, where KGs are often linked to third-party data, reuse standard ontologies and identifiers, and are encoded as RDF, JSON, or other formats. On the right-hand side of the figure, KG maintenance (stage C) is prompted by source updates from stage B, and requirements, audits, and assessments from stage D. To further increase their completeness, correctness, and utility, KGs are refined by completion tasks such as link prediction [90] and error detection and correction tasks, etc [91]. At stage D, there are a selection of use cases for KGs alongside other forms of AI. KGs are used as knowledge bases to query and reason upon, for instance in search [92], question answering [93, 94], and retrieval-augmented generation [3, 95]. Information can be obtained from a graph through deductive (e.g., logical rules) and inductive methods (e.g., as continuous graph embeddings) [1]. Both methods need to be transparent and accountable to the user [96, 97] to be trustworthy and compliant with laws.

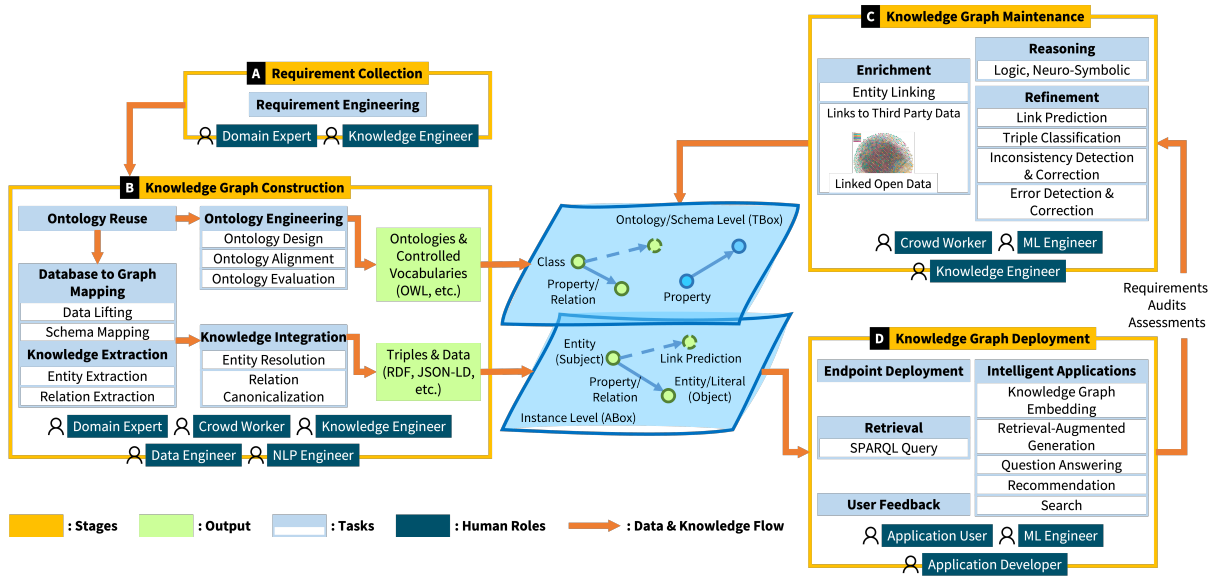


Fig. 1. The knowledge graph lifecycle today.

3. Methodology

To address our four research questions, we employed a mixed methodology of systematic review and interview study. The systematic review involved collecting and analyzing literature on explainable AI in the context of knowledge engineering to gain insight into its current development. The interview study allowed us to directly explore the role of explainable AI in broader contexts, understand the needs of knowledge engineering and KG stakeholders for explanations, identify potential gaps and challenges in this field, and provide valuable insights for further research.

3.1. Literature Review

3.1.1. The PRISMA-guided Review

Following the discussion of the lifecycle, we carried out a PRISMA [98] literature review on databases including ACM Digital Library, IEEEExplore, ScienceDirect, arXiv, SpringerLink, and Google Scholar. We searched for queries combining, on the one side, keywords related to trustworthy (mainly transparent and explainable/interpretable) and, on the other side, keywords related to KG construction tasks, as shown in Table 1. The search initially encompasses all keywords related to KG construction tasks, as depicted in Figure 1. We conducted a prototype search by examining the top 20 results generated by these keyword patterns. Subsequently, we eliminate

Domain	Knowledge Graph Construction	Transparency AI
Keywords	knowledge graph construct*, knowledge graph develop*, knowledge graph complet*, knowledge graph refine*, knowledge graph reasoning, knowledge graph inference, knowledge engineering, named entity recognition, extract entit*, relation extract*, entity linking, entity matching, entity resolution, entity alignment, link prediction	transparent, transparency, interpretable, interpretability, explainable, explainability

Table 1

Keywords for the literature search query. Keywords from two groups were combined for query construction. '*' represents wild characters that can match any word suffix in the search.

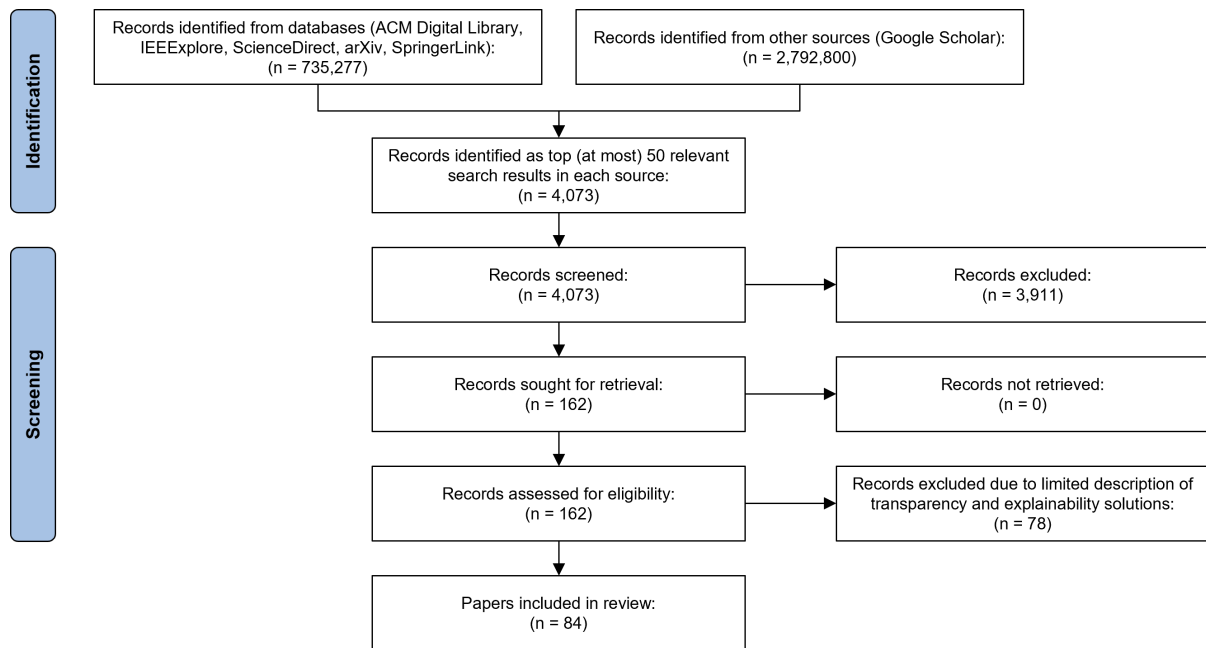


Fig. 2. The PRISMA flow diagram for systematic review.

keywords associated with tasks that do not yield hits within the top 20 results, thereby streamlining the review process. The search took place from October to December 2022 and resulted in more than 735K hits. We then took the top 50 hits per query, which led to around four thousand papers with duplicates⁶. The workflow of paper selection is shown in Fig. 2. We assessed relevance based on titles, abstracts, and keywords first, and in a second step, reviewed the text of the paper to select only those papers that proposed a solution to transparent and explainable KG construction, either as a whole process or for individual tasks. We discarded papers that only mentioned transparency and related concepts rather than putting forward a solution. The final corpus consisted of 84 papers. The papers were all published in the past ten years, which was to be expected given the term "knowledge graph" was coined in 2012 and is in line with other recent knowledge-graph surveys [12, 63].

3.1.2. Use Case Analysis

In addition to reviewing the existing work categorized in Section 4.1, we also analyze the capability of explainable techniques for constructing knowledge graphs through an examination of various use cases inspired by [99]. Specifically, we have identified and presented four use cases along with their objectives in Table 2 and used them as a lens to analyze the papers.

⁶The six platforms where we performed the search supported different query affordances. This means that in some cases, it was possible to build complex queries with multiple keyword options, whereas in others, we had to use separate queries to achieve the same results. We took the first 50 hits for each search query.

Use case	Intentions
Model selection and building	Help knowledge engineers understand the characteristics of ML models thus select the fitted model and build the pipeline
Model debugging	Detect errors may happen in the process and help model users avoid or fix the error
Understanding performance and contributing factors	Give explanations to the predictions and analysis to the contributing factors of the final results
Managing updates	Help knowledge engineers understand how the pipeline will change as data updates and help improve the results

Table 2

Summary of use cases of XAI methods in knowledge graph construction process and their related objectives.

Use Case 1: ML Model Selection and Building When ML is incorporated into knowledge engineering, ML and knowledge engineers must select the proper models and build them. To help users evaluate and select suitable ML models, explanations should reveal the characteristics and limitations of the model, potential risks associated with its use, and its specialization for data or domains. In particular, they should address questions such as the model's capabilities, strengths and weaknesses, and data fitting. It is also important to determine if the model exhibits bias toward specific groups of data sources.

Use Case 2: ML Model Debugging One of the purposes of providing explanations for ML models is to facilitate debugging by allowing knowledge engineers to identify inaccuracies and flawed predictions and providing them with actionable information to correct them.

Use Case 3: Understanding Performance and Contributing Factors To ensure a thorough comprehension of performance, explainable knowledge graph construction pipelines should include the following elements:

- A clear understanding of the inference/reasoning process, which can be represented as rules, paths, etc.
- Identification and highlighting of the factors, important features, and supporting evidence that contribute to the final predictions.
- Provision of counterfactual interpretation through perturbation/permutation.

Use Case 4: Managing Updates Explainability is crucial for knowledge graph maintenance. When updates occur in data sources and contextual information, the knowledge graph can be updated by rerunning the construction pipeline, executing update or modification models, and so on.

To validate the use cases, we compared the use cases derived from the literature review to the ones collected from the interview study and found that the use cases derived from the literature review are mostly reflected through the interview study, and the latter also provide new ones, which we further discussed in Section 4.2.

After identifying use cases, we conducted further investigations into the capabilities of existing works with respect to these use cases. There are two main aspects to consider for this purpose. Firstly, we need to determine whether the reviewed methods have been applied in real-world scenarios of the given use case or could be adapted to suit them. Additionally, we need to consider whether the models have been trained and tested on real-world data. In the domain of knowledge graph construction, benchmarks and datasets are usually close to real-world KGs, such as Wikidata, DBpedia, and Freebase. The second aspect to consider is whether the explanations provided are understandable and satisfactory to the intended audience for the given use case. This can be determined if the work has done comprehensive evaluations that include metrics and human evaluations. Thus, we will evaluate the capabilities of the existing methods based on the following criteria:

- ✕: It is not clear if the method is applicable to the given use case.
- ☆: The method has potential for the given use case.
- ★: The method has been applied to the use case but has not yet been integrated into toolkits or applications in real-world scenarios. Additionally, the explanations provided by the method have not been evaluated through user studies or evaluations.
- ★★: The method has been integrated into toolkits in real-world scenarios. Furthermore, the explanations provided by the method have been tested through real-world studies with the target audience.

3.2. Interview Study

Besides the literature review, we conducted semi-structured interviews with the objective to (1) acquire a basic understanding of the current status of knowledge engineering tools, including transparency issues and obstacles, (2) figure out gaps between existing solutions and practical knowledge engineering scenarios, (3) collect practical requirements for explainable capabilities, and (4) capture insights to design automated explainable knowledge engineering pipelines. Table 3 lists all participants and their background information⁷. In total, we interviewed 13 researchers and knowledge engineers from August to November 2023. All participants were recruited via contact lists of research events, a hackathon, and mailing lists hosted by W3C⁸. We kept a balanced background of the participants on gender (6 females and 7 males) and sectors (7 from universities, 2 from research institutes, and 4 from companies), while achieving diverse coverage in experience, domain, and tasks. Each interview lasted 35 to 50 minutes via an online video call, which involved the authors and the participants. The ethical clearance was granted from the Research Ethics Office of King’s College London with ethics registration confirmation reference number MRSP-22/23-34456.

ID	Job role	Experience	Domain	Tasks
A	Researcher	2	Culture	Ontology engineering
B	Researcher	9	Legal, finance, culture	Link prediction, knowledge extraction, entity resolution
C	Researcher	1.5	Culture	Ontology engineering
D	Knowledge engineer	23	Medicine, scholarly, industry, finance, etc.	Ontology engineering, tool development
E	Researcher	3	Social science	Knowledge extraction, entity resolution
F	Researcher	11	Industry, environment, tourism	Ontology engineering
G	Researcher	9	Public knowledge graphs	Knowledge extraction & completion, entity resolution, ontology matching
H	Researcher	17	IoT, medicine, insurance, tourism, etc.	Ontology engineering, data transformation
I	Knowledge engineer	10	Customer data, public knowledge graphs, geography	Knowledge extraction & enrichment, data transformation
J	Researcher	10	Scholarly, cross-domain	Subject indexing, ontology engineering
K	Knowledge manager	3	Cross-domain	Ontology design, communication
L	Researcher	20	Mobility, manufacturing	Ontology engineering, knowledge extraction, data transformation
M	Researcher	11	Biology, property, medicine, legal, energy, history	Knowledge extraction, knowledge completion

Table 3

Background information about interview study participants includes the job role, experience working with knowledge graphs and knowledge engineering in years, domains of knowledge graphs, and tasks involved in the knowledge graph lifecycle.

3.2.1. Interview Questions

Table 4 presents all the interview questions organized by topics and the order in which they were asked. The questions addressed various topics, including the understanding level of tools and methods, degree of automation, data provenance and lineage, trust, evaluation and human intervention, explainability, and associated risks. The design of interview questions incorporated multiple factors, drawing from previous interview studies on explainable AI in other fields [55, 56], taxonomies and surveys of transparency and explainability [51], and the Explanation Ontology [100] to ensure comprehensiveness. We adapted these trustworthy factors to the context of knowledge graph construction. Firstly, we asked questions about the research background, including experience and domain, to acquire demographic information. Next, we asked about the participants’ experience and understanding of the tools they use. This foundation allowed us to assess the extent to which transparency is an issue and its impact on their practical work. Given the importance of data provenance as a dimension of transparency [101], we include questions specifically about this. To examine the human role in knowledge engineering and gain insight into human factors, we asked questions related to the evaluation of results and how humans interact with the pipeline, providing oversight and intervention. Inspired by [55], we designed questions about explanation scenarios and use cases. These questions delved into scenarios where participants explain results or models to their stakeholders, seeking to identify explainability concerns, challenges, and requirements. Finally, we addressed risk concerns that might

⁷As a knowledge manager, participant *K* is responsible for designing and consulting on taxonomies and ontologies, as well as communicating and educating about knowledge graphs through presentations, webinars, and writing.

⁸semantic-web@w3.org, public-lod@w3.org

1 arise if transparency and explainability are provided with current tools, ensuring a comprehensive understanding of 1
2 potential issues. 2

3 Furthermore, by selecting examples from the previous literature review, we designed XAI examples and facili- 3
4 tated discussions on their usefulness, faithfulness, and acceptance. This approach directly connects stakeholders in 4
5 the context of knowledge engineering with existing literature methods, highlighting the pros and cons of current 5
6 explainable solutions and the gaps between these solutions and practical needs, given the limited application of 6
7 existing XAI approaches in real-world knowledge engineering scenarios. The XAI examples were directly selected 7
8 from the reviewed papers. We first identified papers that provided examples of explanations, such as visualizations 8
9 of attention weights, graph paths, and tables of reasoning rules. We then randomly selected two papers per task 9
10 as examples for participants to discuss. During XAI example discussions, participants were first asked to select 10
11 one (or two, if time permitted) task that they were familiar with. We then provided two examples of two explain- 11
12 able approaches to the selected task. Each example was presented on a slide, consisting of the input, output, and 12
13 explanations as provided by the original publication. Table 7 lists the examples we selected, along with their rep- 13
14 resentations and citations. After reviewing the examples, participants discussed the usefulness and acceptance of 14
15 the explanations, such as whether they found the explanations helpful and whether they would accept them in their 15
16 work scenarios or expose them to stakeholders, such as domain experts and users. Moreover, they were encouraged 16
17 to identify defects in the explanations and suggest improvements or alternative solutions to make the explanations 17
18 more acceptable. During this process, participants were free to ask questions about the provided examples, and we 18
19 responded based on the original publication. 19

20 3.2.2. Coding and Analysis 20

21 The interviews were recorded using Microsoft Teams and transcribed with its automatic transcription services. 21
22 The transcripts were then further cleaned and edited by the authors to remove repeated words, pauses, filler words, 22
23 and to recover errors such as software names and abbreviations. The edited transcripts were coded into keywords 23
24 and patterns, consisting of phrases and sentences. We employed three levels of coding strategies for different types 24
25 of questions. First, for questions related to background information, domain and tasks, and status, we used in vivo 25
26 coding, extracting the exact words from the transcripts. For questions on data provenance and lineage, evaluation 26
27 and human intervention, explanation scenarios, and requirements, we extracted the phrases and identified patterns 27
28 such as operations, methods, and examples. Finally, for questions on understanding, XAI example discussions, and 28
29 risks, we extracted patterns such as comments and suggestions, and coded the attitudes and beliefs towards the 29
30 explainable examples. To analyze the coded data, we grouped identical and similar content into clusters of thoughts 30
31 and insights, and counted the occurrence of each cluster. We also highlighted quotes to provide important supporting 31
32 evidence, insights, and original ideas. 32
33 33
34 34

35 4. Findings 35

36 4.1. The Status of Explainable Automated Knowledge Graph Construction 36

37 4.1.1. The State-of-the-Art XKGC Methods 37

38 We classified the papers reviewed with respect to the KG construction tasks they addressed and their approach 38
39 to explainability, starting with categories widely used in the literature. For explainability, we started with what is 39
40 explained: *local* (data point) vs. *global* (outcome); and when: *post-hoc* (after prediction) vs. *self-explaining* (while 40
41 predicting). We then added another layer for post-hoc methods, splitting the methods into two subgroups: *model-* 41
42 *specific* (specific to one or a group of models) and *model-agnostic* (can be applied to any model). 42
43 43
44 44

45 The results are presented in Figure 3 and visualized using a Sankey diagram in Figure 4. At a glance, the papers 45
46 do not cover the entire KG lifecycle. Most papers are concerned with knowledge acquisition via entity extraction (as 46
47 a source of classes and instances in KGs) and relation extraction (as a source of property classes, but more impor- 47
48 tantly connecting entities to each other through properties), or with curation and maintenance via entity resolution 48
49 (consolidating the data that refers to the same entities) and link prediction (suggesting missing or emerging facts). 49
50 Besides the four core tasks in the bottom half of the figure, we found one paper dealing with the evolution of the KG 50
51 51

Topics	Questions
Background information	What is your job title?
	How long have you been working on knowledge engineering and knowledge graphs?
Domain & Tasks	Can you give a brief description of your work with knowledge graph construction, including:
	<ul style="list-style-type: none"> • an introduction to the knowledge graphs, their types and domains, • tasks in which you have been engaged, such as knowledge extraction or completion?
Status	Could you please briefly describe the tools and methods that you use for the tasks you mentioned?
	<ul style="list-style-type: none"> • Are they fully automated or incorporate human efforts, e.g., human-in-the-loop? • Are they explainable or transparent? And why? • How do you perform the model selection?
	Do you understand, or do you need to understand how the automated components work in detail?
Understanding	<ul style="list-style-type: none"> • What are the obstacles to understand the performance of the component or the results it generated? • If not understood, will the opaqueness of the toolkits impact your work?
	Do you know where the data comes from?
Data Provenance & Lineage	Do you keep track of all operations that have been carried out?
	How do you keep track of data provenance and lineage?
Evaluation & Human Intervention	How do you verify or evaluate the results generated by the automated components or the pipeline?
	Are there any mechanisms to help you?
	If you could verify the results, is there any way that you can correct or modify them?
	What kind of intervention do you take? Explain when and how you perform the oversight.
Explanation	Do you explain to another person how the automated components work or the generated results?
	To whom do you explain the components or results?
	What type of content do you explain?
	How do you explain the results? Do you adopt any methods to help you deliver the explanation?
Use Case	Do you encounter any challenges in this process?
	In what scenarios would you need the pipeline to give you an explanation?
XAI Example Discussion	Please select one of the following tasks, we will provide explainable examples and we can discuss them: (1) entity extraction, (2) relation extraction, (3) entity linking, (4) link prediction, (5) inconsistency detection
	After answering and thinking about the above questions, how would you envision a solution?
Requirements	What kind of information do you hope to be provided by explainable pipelines?
	What is your preferred form of explanation?
	How will the explanations help your work?
Risk	What is your concern about the risk of explainability or transparency?

Table 4

The list of interview questions. The XAI example discussions are accompanied by slides introducing and showing explanations, and the following questions in this part are mostly intrigued by the responses of participants.

schema or ontology [180] and another one about detecting and explaining inconsistency in KGs [181]. We note that link prediction was by far the most popular task, and that a majority of papers dealt with curation and maintenance rather than building a KG for a particular purpose. This is somewhat concerning, as many applications of KGs are in enterprise contexts [6], where the first step is to build a computational representation of the enterprise's data, which is stored across various systems and modalities. We argue that for the tasks not included in the review, there are several potential reasons why almost no papers were found. Many of these tasks still rely heavily on manual work and human oversight and have not yet been automated, as we will later verify based on interview results. This includes tasks such as ontology reuse and ontology design. Additionally, there are tasks where automation, such as the use of LLMs, has been employed, like ontology alignment [185] and data lifting from databases, but explanations have not been considered.

A second high-level observation is the balanced split in the chosen format for explanations. Methods based on input and generated features use attention weights [124, 131], words [117, 118], attributes [102], etc. to generate explanations, which can be numerical, textual, or visual. By contrast, methods based on human-understandable background knowledge provide explanations in the format like logical rules [164], reasoning paths [160], and structured contextual information [120] as explanations. Given that we are interested in explanations that are accessible to knowledge engineers and subject domain experts, it would be interesting to evaluate if their familiarity with knowledge representation and/or the subject domain impacts how useful knowledge-based explanations are compared to

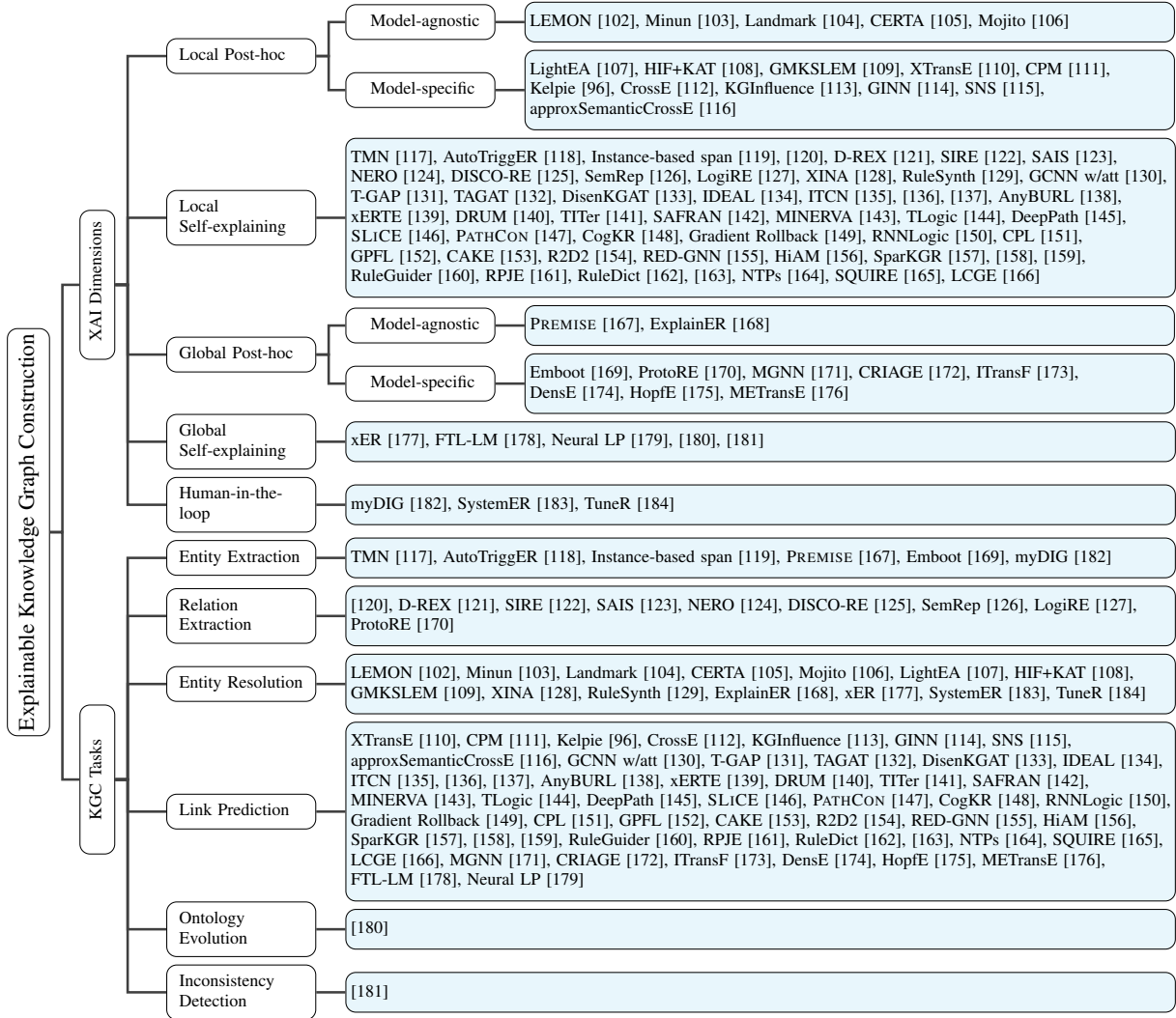


Fig. 3. Taxonomy of Explainable Knowledge Graph Construction.

feature-based ones, which sometimes require an understanding of machine learning. At the same time, explanations are generated in a different way for each of the four core KG construction tasks in the bottom half of the figure.

Entity Extraction For entity extraction, explanations often leverage contextual cues such as triggers [117, 118] and patterns of words [167], utilizing attention mechanism [186] and saliency map techniques. One notable work is myDIG [187], a human-in-the-loop system that compiles sophisticated rules written by domain experts into SpaCy rules for backend execution. This reduces the barrier for domain experts to interact with the machine and minimizes training effort. Additionally, myDIG records extraction provenance, allowing users to explore the downstream effects of their specifications. Another type of explanation used for entity extraction is example-based explanations, which rely on training instances [188]. In Ouchi et al., similarities between pairs of candidate(s) and the training instances are computed, with the term having the highest derived label probability being returned [119].

Relation Extraction For relation extraction, explanations frequently employ contextual information from the input, such as words and sentences, similar to entity extraction. The attention mechanism is a prominent principle among relation extraction methods, with 4 out of 9 studies using attention weights and their associated input context to generate explanations. For instance, NERO uses word-level attention to calculate matching scores between sentences

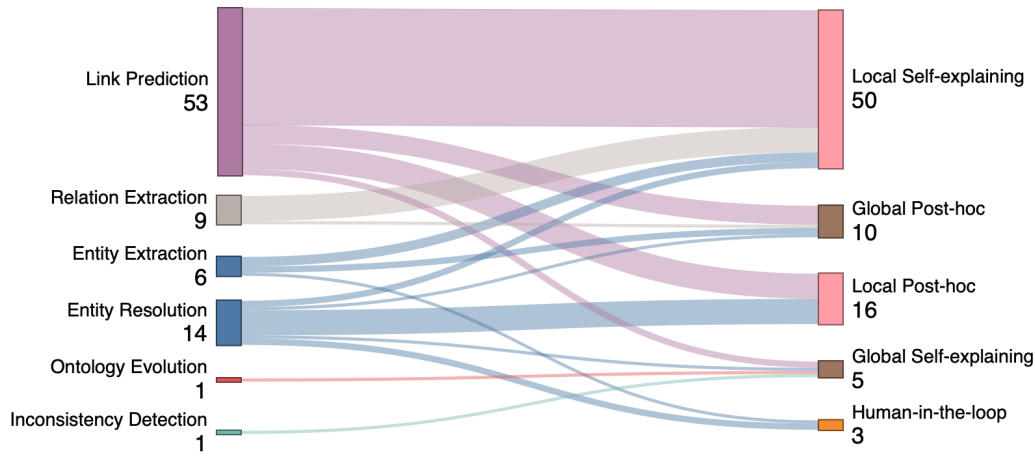


Fig. 4. Sankey diagram of methods categorization. The left column represents KGC tasks and the right column represents XAI taxonomy. The total number of each categorization is labeled under their names.

and generated rule patterns, where attention weights represent word importance for constructing attention-pooled rule/sentence representations. SIRE employs the attention mechanism in both the evidence selector [189] to identify supporting evidence and in the logical reasoning module [186]. In light of research questioning the validity of attention as faithful explanations, Shahbazi et al. adopted a mixed explanation mechanism extended by saliency, etc [120]. Another prevalent type of explanation in relation extraction involves relation learning/logic rules. Beyond NERO, LogiRE integrates a rule generator and a relation extractor, optimizing these modules using the expectation-maximization algorithm for document-level relation extraction. Diverging from text-based explanations, ProtoRE learns prototypes for each relation from contextual information, exploring the intrinsic semantics of relations and visualizing them as geometric explanations [170]. Under optimal conditions, prototypes are unit vectors uniformly dispersed on the surface of a unit ball, with datasets clustered around each prototype vector.

Entity Resolution There are two primary types of explanations for entity resolution: entity matching (EM) rules [108, 129, 183, 184] and (ranked) attributes of the entity pair with relevant scores [102, 104–106, 168]. EM rules, represented in forms such as disjunctive normal form and general boolean formula, are commonly used in EM systems to enhance interpretability [190]. For automatic EM rule-based models, Yao et al. proposed a framework consisting of Heterogeneous Information Fusion for learning feature representation from unlabeled data and Key Attribute Tree for interpretable EM decision making [108]. This framework translates decision trees into EM rules, making explanations more accessible to domain experts. RULESYNTH, proposed by Singh et al., formulates the rule discovery problem in entity matching as a program synthesis problem. They adopted a more concise and interpretable form of General Boolean Formula to represent EM rules and proposed a novel rule synthesis algorithm. In contrast to EM rules, using attributes and their relevant scores as explanations focuses on uncovering the contribution and importance of each attribute or combinations of attribute sets in the decision-making process of entity matching. Most works employing this representation adopt perturbation-based methods. By applying LIME perturbation to the entity resolution problem [42], these methods use perturbations such as dropping words and making entity pairs less similar to analyze differences and calculate predefined importance scores. Additionally, they explore the presence of input attributes and attempt to add them to make non-matching pairs more similar.

Link Prediction Most explainable link prediction methods leverage the topology and reasoning capabilities of knowledge graphs (KG). Rule- and path-based methods have become the predominant forms of explanations, achieved through various approaches such as random walk-based methods [138, 144, 178], reinforcement learning agents [143, 145], and perturbation-based methods [96, 172]. A significant body of work utilizes reinforcement learning (RL) for reasoning over knowledge graphs and searching for paths to explain link prediction results [136, 141, 143, 145, 151, 154, 157, 160]. These models typically comprise knowledge graph environments and policy network agents. The knowledge graph environment transitions elements within the graphs (e.g., entities,

relations, queries) into RL agent elements, where states are usually entities (in practical terms, embeddings) and queries (subject entities and relations); actions are typically outgoing edges/relations; transitions map current entities and their outgoing edges to their neighboring nodes; and rewards are heuristic indicators, awarding 1 when the agent reaches the correct target entities. Policy networks then maximize the expected reward to perform path finding. Variations exist in environment transitions, rewards, and the parameterization of the policy function. For example, R2D2 [154] and RuleGuider [160] employ multi-agent architectures. R2D2 uses two agents, with one arguing the fact is true and the other arguing it is false, feeding their arguments into a judge network. RuleGuider uses a relation agent and an entity agent that interact to generate paths fed into a rule miner. Perturbation-based methods are also applied in link prediction, similar to those used in entity resolution. CRIAGE [172] introduces graph perturbation by removing a neighboring link from the target fact to assess the influence of the fact and by adding a new, fake fact to evaluate model robustness and sensitivity. Another prevalent method in explainable link prediction models is the attention mechanism, used in 16 out of 53 total link prediction works. For instance, XTransE employs attention values on items to reveal the relevance between different property-value pairs and the current prediction, which are then ranked to identify the most relevant triples [110]. In xERTE, Han et al. propose a temporal relational graph attention layer that calculates query-dependent attention scores for each edge [139]. These scores propagate to each node's prior neighbors, pruning the inference graph using edge contribution scores. The pruned graph, with node attention scores and edge contribution scores, is used to produce the explanations.

Human-in-the-loop There are very few papers considering human inputs or oversight, which are critical in trustworthy AI frameworks and guidance [191]. In the few cases of human-in-the-loop systems, human input often involves the provision or revision of rules for tasks such as entity extraction [182] and entity resolution [183, 184]. In myDIG [182], a GUI-based rule specification system is provided for domain experts to input expressive entity extraction rule sets without programming. SystemER [183], which adopts an active learning methodology, learns explainable entity resolution logical rules and offers functionalities for domain experts, both with and without programming backgrounds, to verify and customize the learned models in feature engineering to ensure extensibility. For generating entity resolution rules, TuneR [184] involves developers (i.e., coders, scientists, and domain experts) in tuning rule sets by defining the contribution of optimization metrics. The framework defines interpretability-related metrics as the preference between the number of rules in the rule set and their overlap. All three approaches use an ensemble of rules to achieve high precision. Several factors influence the success of these human-in-the-loop approaches, some of which have been considered in these three systems. One critical factor is balancing the minimization of training with the extent of human intervention. More human intervention can reduce training efforts, which require feeding more data and extending training time. Conversely, increased training efforts can reduce human intervention, thereby minimizing unnecessary human labor and avoiding time-consuming and error-prone trial-and-error processes. Another factor is the degree of operational freedom given to users. The relationship between the complexity of functions and the freedom of operations provided to users affects the time required to educate them. The design of functions should enable users to maximize their input to produce high-quality work while minimizing the time needed to familiarize themselves with the tool. Providing too few intervention options might hinder users from fully expressing the correct input, thereby increasing human effort. These factors are crucial when designing human-in-the-loop systems, and more user studies, especially for knowledge engineers and knowledge graph stakeholders, are needed to explore them further.

Evaluation of explanations We also collected and analyzed the evaluation of explanations. A primary observation is that most XAI approaches have not been thoroughly and/or comprehensively evaluated. The majority of methods (58 out of 84) do not perform any evaluation on explanations or only use anecdotal evidence by visualizing and commenting on a limited number of cases of explaining outcomes intuitively. There are efforts to design metrics to evaluate explanations. 17 works adopted metrics to evaluate their explanations, and most of them are task-dependent. Shahbazi et al., [120] created a ground-truth explanation set and computed the Kendall Tau correlations for the sentence importance scores for the annotated test set. approxSemanticCrossE [116] proposed explanation evaluation metrics target the link prediction tasks, which calculate the ratio of triples for which the model can generate explanations (recall) and the number of explanations, on average, for each prediction (average support). In gradient rollback [149], Lawrence et al. adopted the "RemOve And Retrain (ROAR)" [192] evaluation paradigm to evaluate the faithfulness of the explanations.

Evaluation tasks	Methods	Number of participants	Background
Comparing model-generated and human-provided explanations	AutoTriggER [118]	/	crowd-workers
	D-REX [121]	3	crowd-workers
Judging the relevance and correctness of explanations (examples)	AutoTriggER [118]	See above	
	Emboot [169]	2	domain experts
	xERTE [139]	53	*
	DRUM [140]	2	CS students
Comparing explanations generated by different models	D-REX [121]	See above	
	RuleSynth [129]	27	CS researchers
	RuleGuider [160]	/	crowd-workers
	DRUM [140]	See above	
Survey with questions measuring the quality (usability, reliability, trust, etc.) of explanations	SQUIRE [165]	/	authors
	Kelpie [96]	44	/
Evaluating the accuracy or precision of user predictions with or without explanations	SQUIRE [165]	See above	
	R2D2 [154]	44	/
	[137]	/	domain experts

Table 5

Works that use human evaluation to analyze explanations. “*” indicates that no group label is provided, but other detailed background information of participants is reported. ‘/’ means ‘not reported’ in the paper.

12 studies use human evaluation, detailed in Table 5. We identified 5 types of evaluation tasks commonly adopted in these studies. The most frequent tasks involve asking participants to compare model-generated explanations with those from baseline models and to judge the relevance and correctness of a set of examples. Various metrics are used in human evaluations. One approach is to have participants rate the usability, reliability, and trust of explanations in a survey. A notable example in this group is SQUIRE [165], which annotates BIMR-based interpretability scores [193] for paths generated by their models and baseline models. Another group of methods measures the accuracy or precision of user predictions with or without provided explanations. The backgrounds of human evaluators are varied, including domain experts, such as e-commerce experts in [137] and linguists in Emboot [169], people with technical backgrounds, and laypeople such as crowdsourcing.

From the above observations, we identified several issues with the evaluation methods. First, reporting a limited number of examples selected based on the researchers’ intuition can be biased and not sufficient for robust verification [57, 194]. Since not all results have satisfied explanations generated, another issue is that the ratio of results for which the model can generate satisfied explanations is not commonly reported. In our interview study, we found it to be a crucial factor that might influence the user’s trust in the XAI models.

4.1.2. Use Cases and Capabilities Measurement

The capability of various explainable techniques for each use case is shown in Table 6. In general, the reviewed literature indicates that global post-hoc methods, especially model-agnostic ones, have the potential to address all use cases. Local post-hoc methods have also demonstrated similar potential across all use cases. Although no global self-explaining methods were identified for the first two use cases, this does not imply that these methods lack potential for model selection, construction, and debugging. Instead, they are suitable for providing model analysis due to their global assessment capabilities. Among the use cases, three areas, besides understanding performance and contributing factors, have received less attention and research. This could pose challenges when integrating developed methods into real-world applications, making it essential to address these gaps.

Use Case 1: ML Model Selection and Building Most model-agnostic methods, such as explainers designed for any knowledge graph embedding models and some model-specific methods, have the advantage of providing explanations across different models and facilitating comparison. While some of the reviewed works have demonstrated their applicability in this use case, most have not emphasized addressing concerns related to model selection and comparison. A notable example that covers this use case is ExplainER [168], which offers a mechanism for model analysis. The analysis engine of ExplainER comprises multiple explanation tools (LIME [42], Anchors [195], BRL

Use case	Local Post-hoc	Local Self-explaining	Global Post-hoc	Global Self-explaining	Methods
Model Selecting and Building	★	☆	★	×	ExplainER [168], CPM [111], Kelpie [96]
Model Debugging	☆	☆	★	×	LEMON [102], ExpalinER [168], D-REX [121], Instance-based [119], [182], TuneR [184], CRIAGE [172], Kelpie [96], SparkGR [157], GCNN w/att [130], MINERVA [143]
Understanding Performance and Contributing Factors	★★	★★	★★	★★	All papers
Managing Updates	☆	★	★	★	ExplainER [168], CPM [111], TuneR [184], Abstraction [181], TLogic [144], Emboot [169], SystemER [183], RNNLogic [150], [137] Neural LP [179], FTL-LM [178], ITCN [135], CPL [151], SQUIRE [165], MGNN [171], DRUM [140], ProtoRE [170], CRIAGE [172], TITer [141], PATHCON [147], METransE [176], RED-GNN [155], [180]

Table 6

Capabilities of XAI methods in knowledge graph construction. Symbols are referenced from Section 3.1.2: ×: applicability to the given use case is unclear; ☆: method shows potential for the given use case; ★: method has been applied to the use case but is not yet integrated into toolkits or real-world applications. Explanations provided by the method have not been evaluated through user studies or any other evaluation methods; ★★: method is integrated into toolkits in real-world scenarios, and its explanations have been tested through real-world studies with the target audience.

[196], and Skater [197]) that are independent of any entity resolution models. For link prediction, explainable methods such as CPM [111] and Kelpie [96] can be used with any embedding-based link prediction models, allowing for comparison across different embedding models. The main gap for current models in this use case is not solely related to model design and architecture, but also to better documentation. One potential solution is to document an interactive model card [198] that lists all the necessary information regarding explainability. For instance, for explainable link prediction models, this could include the ratio of faithful and correct explanations generated for each embedding model and a comparison of generated explanations for the same input.

Use Case 2: ML Model Debugging Some works provide analyses of errors. For example, the instance-based explainable method performed error analysis using relevant examples to identify factors causing model confusion [119]. ExplainER visualized representative explanations to highlight where the model fails [168]. D-REX conducted error analysis on explanations alongside model predictions, further revealing the model’s error detection capabilities [121]. Pezeshkpour et al. demonstrated the potential application of CRIAGE for automated detection of erroneous triples in knowledge graphs. Their approach focused on identifying triples with the least influence on the model’s prediction of the training data [172]. Similarly, Rossi et al. highlighted the ability of Kelpie to uncover bias and imbalance in data, enabling researchers to correct it. However, although these works provided analyses of errors, most did not offer actionable steps for rectifying the identified issues. This could be achieved by providing options to adjust parameters, model architectures, and leverage external sources such as human knowledge. Human-in-the-loop methods exemplify approaches for correcting errors and improving model output, such as manually correcting rules by domain experts in rule-based explainable systems [182, 184]. One approach following this line is to offer local actionable information, such as suggestions for correcting predictions directly. A future direction in designing local explainable methods would be to help users identify error cases and enable corrections at the data point level.

Use Case 3: Understanding Performance and Contributing Factors The majority of the reviewed work performed well across various tasks. As detailed in Section 4.1.1, a range of representations are employed to understand the inner workings of models and the factors contributing to their outputs. For knowledge extraction tasks, such as entity and relation extraction, models provided supporting evidence from the source data (e.g., text) to aid in predictions [117, 118]. Similarly, for knowledge integration tasks like entity linking, attributes of entities were selected through mechanisms such as matching or non-matching votes, as demonstrated in [102, 104]. Explainable link prediction models offered rules [129, 138, 160] and paths [143, 145, 151, 154] to illustrate the reasoning process, as well as subgraphs [148] to measure the influence of nodes and edges. Notably, rule-based methods are prevalent across all tasks due to their concise and straightforward representation and their ability to generalize to new data.

Use Case 4: Managing Updates Global explainable methods such as rule-based methods [144, 184] can potentially express the model evolution through modifications in their global explanations. Similarly, visualization-based explanations [170, 176], where users can compare different versions of visualizations, can also provide valuable insights when managing updates to knowledge graphs. Models that provide local explanations, such as inductive models [140, 141, 147] and perturbation-based models [172] could track differences for specific instances or groups of instances. Very few of the models directly implemented this capability, but most of them could be potentially extended to support this use case. For rule-based explainable methods, a straightforward way to manage updates is to use the generalization ability of existing rules and perform inductive reasoning. For instance, TLogic [144] stated that the temporal rules they generate are applicable to any new dataset, as long as the new dataset covers common relations, even in cases where new entities appear. Zhang et al. [137] also emphasized the benefits of transferable rules. Their model could generate reusable rules to accelerate the deployment of a knowledge graph to new tasks or systems. In addition to directly transferring rules to new data, rules can also be updated. For example, RNNLogic [150] used an EM-based algorithm to update rules. Once the explanation rule sets were updated, to gain more insights, the users could compare two sets of rules and see what changes the new data had brought in. Similar strategies can be applied to other explanations, such as the visualization of attention weights and embeddings.

4.2. Explainable Knowledge Graph Construction in Practical Scenarios

We now report on the interviews. We first present the current status of knowledge engineering tools in practical scenarios, focusing on the degree of automation and the level of understanding that knowledge engineers have towards these tools, as well as aspects including data provenance and lineage, evaluation, and human intervention. By addressing a series of sub-questions, we aim to gain a basic understanding of these critical transparency factors from our interview study. This foundation will enable us to delve deeper into identifying the desired properties of explainable tools.

4.2.1. Automation and Understanding

How much human effort is leveraged in the knowledge graph lifecycle? Out of the 13 participants, the majority engage in manual (5 participants) and semi-automatic (5 participants) work, while a minority (3 participants) exclusively utilize automation for the tasks they work on. From the perspective of task execution in ontology engineering, participants predominantly employed manual and/or semi-automatic methodologies. These approaches necessitate extensive communication and collaboration among knowledge engineers, domain experts, and stakeholders, often facilitated through semi-structured interviews. Conversely, for tasks related to knowledge extraction and completion, participants demonstrated a preference for automated tools and models. Methods, which focus on tasks like data transformation that lifts other formats of data into RDF triples through RML mappings⁹ and tools like SPARQL Anything [199], always involved the manual creation of the mappings. One participant assessed the performance of leveraging language models in generating such mappings. Language models in knowledge engineering have enhanced automation due to their user-friendly nature, characterized by simple natural language input and output, which require fewer specialist skills. However, their opacity and tendency to generate hallucinations impact their trustworthiness. When evaluating the outcomes of tools and models, such as the triples generated by knowledge extraction tools, human evaluation is always necessary. This is particularly crucial when dealing with new domains and data, where datasets are lacking.

What is the level of understanding of the tools and models? Participants had varying opinions regarding the impact of the opaqueness of tools and models and the necessity to thoroughly understand them. 6 of them felt that opaqueness did impact their work and emphasized the importance of understanding the models. As participant A highlighted, this importance extends beyond merely explaining why the models produce certain outputs. It also involves helping humans "understand the extent to which these outputs can be trusted"¹⁰ and determining "how they might need to change the way they interact with the model". Participant B also noted that the opaqueness of the tools and models might complicate evaluations, as it becomes challenging to determine how specific inputs influence the

⁹<https://rml.io/specs/rml/>

¹⁰We use double quotes to indicate that the quotes are the original words of the interviewees.

1 outputs. In contrast, the remaining 7 participants were less concerned with transparency issues, feeling that opaque- 1
 2 ness was not a significant problem. They provided several reasons. Two participants stated that only the model’s 2
 3 performance and the final quality of the output knowledge graphs mattered to them. Since they primarily deal with 3
 4 public datasets and transparency and explainability are not within their research scope, they pay less attention to 4
 5 these topics. Three participants mentioned that their tasks are predominantly manual, so transparency and explain- 5
 6 ability are less applicable. Some participants noted that even collaborative projects require some level of explanation 6
 7 for better communications and outcomes between human agents. 7

8 **All 13 participants demonstrated a relatively high level of understanding towards the tools they used,** 8
 9 particularly when these tools were open-sourced and/or came with documentation (e.g., publications, technical doc- 9
 10 uments), or if the tools were self-developed. Two participants mentioned that it is not always necessary to delve 10
 11 into the code level, and sometimes it is challenging to fully comprehend how the models make decisions. However, 11
 12 it is crucial to gain a conceptual understanding of the mechanisms of specific components, their technical limita- 12
 13 tions, and underlying assumptions. 7 participants reported no significant obstacles in understanding the tools and 13
 14 methods. For the remaining 6, the challenges of understanding the tools and models varied. The primary obstacle, 14
 15 mentioned by 3 participants, was the difficulty in understanding errors produced by models, their causes, and how 15
 16 the models arrived at decisions. Participant *B* pointed out that models could be difficult to understand due to their 16
 17 mathematical complexity and insufficient background knowledge in ML and NLP, particularly when it comes to 17
 18 understanding the inner workings of LLMs. Participant *E* noted the difficulty in determining the optimal size of 18
 19 data and model parameters to train models effectively or to transfer them to another domain or input type. Partic- 19
 20 ipant *L* emphasized the challenge of evaluating both the correctness and the completeness of results, noting that 20
 21 both aspects are critically important. Additionally, understanding "what level of quality is good enough for the task" 21
 22 is also challenging. To address these challenges, participant *A* suggested designing models to provide additional 22
 23 outputs that help in understanding the models. This is somewhat achieved by works in the literature review, such 23
 24 as models leveraging attention mechanisms that output attention weights [117], and reinforcement learning models 24
 25 that produce reasoning paths [145] to explain results. For generative models, asking them to generate additional or 25
 26 intermediate outputs, such as reasons for certain outputs could help. However, this is not always technically fea- 26
 27 sible. For example, adjusting embedding models to generate intermediate outputs for model understanding is not 27
 28 as straightforward as with LLMs through chain-of-thought [200]. Another approach proposed by participant *G* for 28
 29 understanding incorrect results is to seek training examples similar to some test data instances. This aligns with 29
 30 existing reviewed example-based explanations [188], such as [119], where similar training instances are returned as 30
 31 explanations for the assigned label of the candidate instance. 31

32 4.2.2. Data Provenance and Lineage 32

33 *Do knowledge engineers know where the data comes from?* **Among the 13 participants, 12 reported knowing 33**
 34 **the sources of their data. Data sources and providers varied, with participants often receiving multi-sourced 34**
 35 **data, depending on the projects they were working on.** As shown in Figure 5 (a), half of the participants used open 35
 36 knowledge graphs and publicly available data. However, this does not mean the data sources are always clear to them 36
 37 and verifiable. Not all participants are aware of how these datasets and knowledge graphs were created. For instance, 37
 38 participant *E* mentioned that they were unaware of how the benchmark datasets were created, but recognized that 38
 39 this data often has significant limitations, such as skewed distributions and incompleteness. Similarly, participant *H* 39
 40 noted that when using external APIs to obtain data, it was unclear where the data originated from. Data can also 40
 41 be collected from domain experts and stakeholders or acquired from partners and collaborators with data sharing 41
 42 policies and platforms. However, participants using data from these sources reported similar challenges: some of 42
 43 them generally did not make extra efforts to understand the origins of the data and how exactly it was selected. 43
 44 Participant *M* also noted the difficulty in assessing the qualifications of annotators when data is manually annotated, 44
 45 as detailed information about their expertise is often unavailable. Participants from industry also collect data from 45
 46 customers, which can be sensitive and requires extra effort for data masking. A minority of participants constructed 46
 47 datasets themselves. 47

48 *How do people keep track of data provenance and lineage?* Among the participants flagging data provenance as 48
 49 essential, 9 actively tracked it in their tasks. Notably, all participants from industry and academia alike recognized 49
 50 the importance of data provenance and lineage and have established methods for documenting these aspects, given 50
 51 51

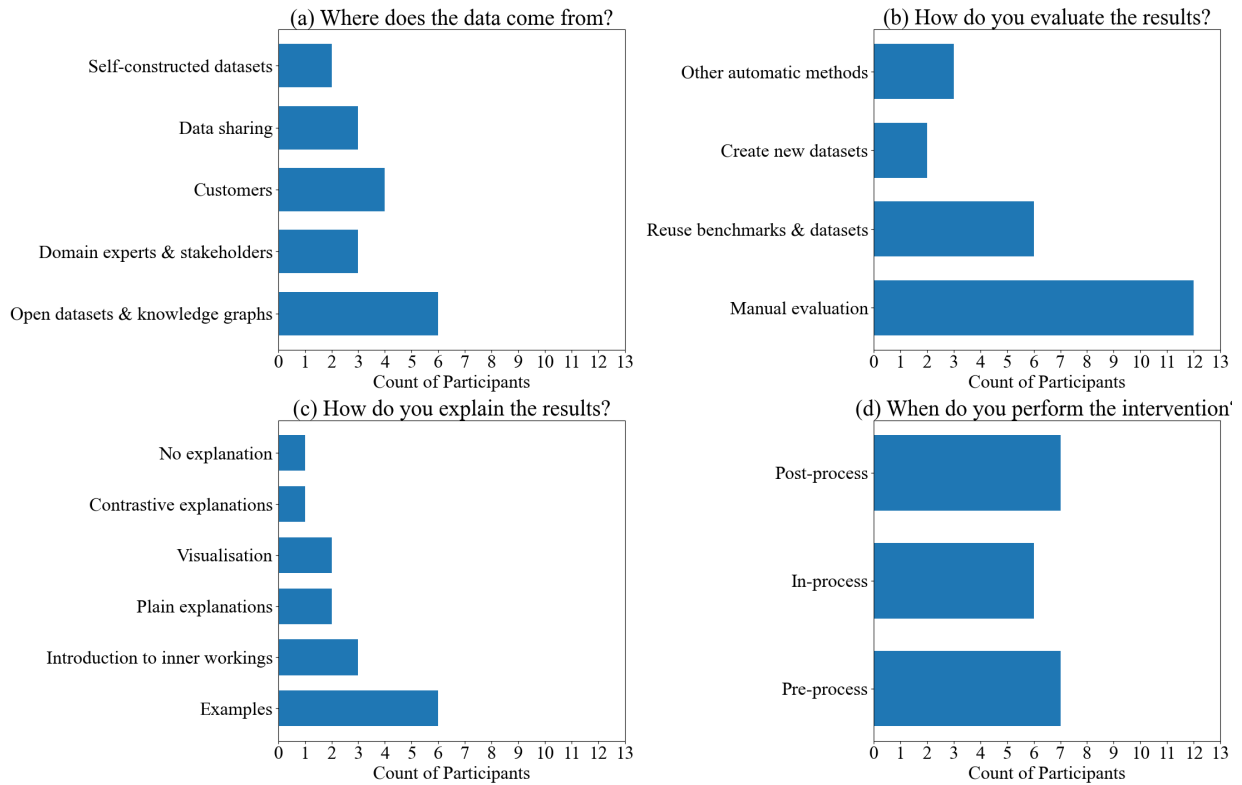


Fig. 5. Distribution of participant responses to four questions: (a) where does the data come from? (b) how do you evaluate the results? (c) how do you explain the results? (d) when do you perform the intervention? The x-axis represents the total number of participants (13), with multiple responses allowed per participant.

that their data primarily comes from partners and customers. The interview revealed a list of (semi-) automatic tools either currently in use or planned for adoption to manage data provenance and lineage by the participants, including PROV Ontology¹¹, RDF Star¹², metadata, OpenRefine¹³, Data Version Control (DVC)¹⁴, data catalogs, NLP Interchange Format (NIF) [201], and blockchain. These tools document a set of details, such as the creation time, involved personnel, operation timelines, algorithms used to create the data, and potentially even the parameterization of these algorithms. Data provenance is tracked at different granular levels, from the model level (e.g., entire ontologies) to the data level (e.g., individual ontology elements). The availability of a wide range of tools offers knowledge engineers flexibility in fitting their specific pipelines. However, challenges and requirements remain. For instance, participant *M* noted the difficulty in determining the extent to which data provenance should be tracked and the level of details required. There is also a preference for using standard representations based on standardized languages and vocabularies. As participant *B* said, an ideal approach would involve "an explicit and standard representation of provenance and lineage directly attached to the produced artifacts", such as "metadata that accompanies the actual knowledge assets". It is also crucial to have an automatic, scalable, and trustworthy method for tracking data provenance and lineage, particularly when dealing with sensitive and frequently updated data. Participants also highlighted that integrating LLMs into the knowledge engineering process introduced new challenges for data provenance, as their extensive training data is often unknown.

¹¹<https://www.w3.org/TR/prov-o/>

¹²<https://w3c.github.io/rdf-star/>

¹³<https://openrefine.org>

¹⁴<https://dvc.org>

4.2.3. Evaluation and Human Intervention

How do knowledge engineers evaluate the results? As shown in Figure 5(b), most participants rely on **human evaluation** to evaluate the outcome of knowledge engineering tasks. The human evaluation methods used are usually qualitative analysis or randomly sampling a subset of data for manual inspection. Depending on the task and the required expertise, the evaluators are usually domain experts and/or the researchers themselves. From our interviews, two participants reported recruiting domain experts for evaluation, one conducted evaluations both by the developers and domain experts, while the remaining 9 evaluated the results themselves. Several reasons contribute to the heavy reliance on human labor for evaluation. Firstly, there is a lack of available datasets and testing platforms. Secondly, existing datasets are often unsuitable for new scenarios. A model that performs well on current datasets may not necessarily perform well on new data, rendering the existing datasets less helpful. Thirdly, metrics used, such as average precision and accuracy, can sometimes be misleading. It is often the case that either the metrics look too good and the results are worse, or the metrics look very bad and the results are better than they appear. There are several additional issues with human evaluation. It is time-consuming, not scalable, and not always feasible. Additionally, randomly evaluating a subset of data might not accurately reflect the actual quality of the generated results. Participant *L* noted, "the main difficulty is to select an actual relevant sample".

Besides manual evaluation, 6 participants also attempt to **reuse benchmarks and datasets** with ground truths to automate the evaluation process. They are aware that benchmarks are incomplete, biased, and skewed, and thus might not always truthfully reflect the models' abilities. Participant *M* from the industry mentioned that there are often gaps between the focus of the benchmarks and the requirements of customers. While benchmarks are more focused on challenging cases and specific domains, customers are often interested in general domain cases, making the benchmarks less relevant to practical scenarios. Participants from academia also face difficulties, as mentioned by *G*, when there are insufficient resources for annotations or re-annotations, forcing them to work with the existing data. Participants also highlighted other methods to automate evaluation. Two participants mentioned using SHACL shapes for validation. Participant *J* stated they are developing a quality management concept and have a machine learning-based algorithm on top of other methods to estimate quality. Additionally, two participants indicated they **construct new datasets** for evaluating their methods. Although there are ongoing efforts to automate evaluation, the interviews revealed a consensus that different extents of human evaluation are always required.

What do people do when they find the results incorrect? Similar to human evaluation, the interviews revealed a consensus that human intervention is essential to compensate for the limitations of machines at various stages and levels of detail in the knowledge graph construction process. We categorized human intervention based on the stages at which it occurs, as shown in Figure 5(d). The first stage is **pre-processing**, where human intervention primarily involves working with the inputs. The most common approaches are data augmentation and cleaning. When working with LLMs, this also involves improving the input prompts. The second stage is **in-process intervention**, where researchers and knowledge engineers adjust the tools and models or specific steps in the process to resolve issues. This could involve fine-tuning and re-training models, debugging code, or adding, removing, and modifying components and steps. Before making these modifications, there is typically a troubleshooting process to identify error patterns, systematic mistakes, and biases. Participant *B* mentioned that when using LLM-based pipelines, disambiguation is always a problem, so they either improve the prompts or add an extra disambiguation step. Another example, provided by participant *I*, is involving humans in identifying incorrect inference rules and then rerunning the models. The third stage involves directly modifying the model output, where humans manually correct a group of generated outputs (if the errors are manageable in size) or add post-hoc filters to exclude problematic results (i.e., **post-processing**). Statistically, 7 participants adopted pre-processing methodologies, 6 engaged in in-process modifications, and 7 employed post-process corrections. This indicates that participants typically adopted multiple types of interventions at various stages. Moreover, the individuals performing the intervention are crucial. It's not just about their availability to check the results, but also about their expertise. As noted by participant *L*, when mistakes are "a mix of technical and domain-specific issues", it can be more challenging to identify and address them.

How do people explain to others their tools and results? Of the 13 participants, 12 have experience explaining models and results to others. Five participants explained their models and results to stakeholders who may not have a technical background, typically domain experts. Eight participants explained their work with ontologists and

1 knowledge engineers, who have a similar technical background, usually project partners and team members. Addi- 1
2 tionally, two participants mentioned producing explanations for educational purposes, targeting university students. 2
3 This indicates that **designing and delivering explanations has become a crucial and challenging task in the** 3
4 **knowledge graph lifecycle**. As highlighted by participant *L*, if the model performance does not meet stakeholder 4
5 expectations and the model is not explainable, it does not foster acceptance or transparency. 5

6 The methods used for explanations are summarized in Figure 5(c). For now, there are no standardized methods 6
7 for explaining the models and outputs in the knowledge graph lifecycle. **We observed that almost no methods** 7
8 **from the literature review are used in participants' daily work scenarios**. We argue that there may be two main 8
9 reasons for this. First, participants may not be aware of these methods and therefore rely on their own intuitive ways 9
10 to explain results when needed. Second, the available methods may not be ready for practical use, and integrating 10
11 existing XAI methods into their workflows is challenging. Only participant *B* mentioned having used one of the 11
12 presented models in the example discussion session [136], finding it generally useful, although not all explanations 12
13 produced by the model were helpful. 13

14 The most frequent method (used by six participants) is to select examples, including corner cases and errors, to 14
15 explain the model's functionality, the relevance between input and output, the difficulties of the problems, and the 15
16 range of the model's abilities. Three participants explained the pipelines and models through lectures and conceptual 16
17 introductions to the technical components, often providing high-level overviews of the algorithms and models. 17
18 The other two participants adopted visualization methods. Participant *A* reported success using visualizations to 18
19 represent embeddings and clusters, which helped "define a clear boundary between technical and intuitive content". 19
20 One participant mentioned using contrastive explanations, such as why the machine made one decision instead of 20
21 another. Two participants did not have a specific method but relied on plain explanations. 21

22 Using the same taxonomy adopted in the literature review in Section 3, we categorized the explanation methods 22
23 collected from the interviews into two categories: contrastive explanations and example-based explanations as lo- 23
24 cal post-hoc methods, and visualization, plain explanations, and introduction to inner workings as global post-hoc 24
25 methods. Our analysis reveals that, out of the 14 responses regarding explanation methods, half of the responses 25
26 are local post-hoc methods, while the other half are global post-hoc methods. Notably, no self-explaining methods 26
27 were reported. In contrast, the literature review indicates that a substantial proportion of explainable methods con- 27
28 sist of local self-explaining (59.5%) and local post-hoc (19%) methods. We posit that several factors contribute to 28
29 this discrepancy. Self-explaining methods are preferred in academia-developed models because researchers often 29
30 work on implementing models from scratch or improving models by adjusting components or integrating additional 30
31 components for better performance. This objective aligns with the design of self-explaining models. Among the 50 31
32 local self-explaining papers reviewed, 37 pertain to link prediction models, which typically incorporate explanation 32
33 mechanisms into their developed models, enhance existing models by making components explainable, or reformu- 33
34 late problems in an interpretable manner. For practitioners, however, implementing self-explaining methods poses 34
35 challenges. Post-hoc explanations of model output offer greater flexibility, allowing practitioners to customize sup- 35
36 porting evidence, visualize this evidence, and adapt explanations into other languages that are more comprehensible 36
37 to their stakeholders. 37

38 Participants reported several challenges. **The most significant challenge when providing explanations is the** 38
39 **gap in background, knowledge, and requirements between knowledge engineers/researchers and their stake-** 39
40 **holders**. This gap makes it difficult to capture the interests and needs of stakeholders and to determine the appro- 40
41 priate level of technical detail for explanations. Reporting too many technical details may disinterest and frustrate 41
42 stakeholders who lack a technical background. As participant *J* mentioned, "they might feel confused and show 42
43 very little interest in what the numbers in the explanations represent." Participant *B* noted that there is often a gap 43
44 in expectations, with knowledge engineers and researchers focusing on "understanding from the technical aspects" 44
45 while stakeholders are more interested in "the practical use case and deployment perspectives." Participant *L* added 45
46 that audiences often have "distorted expectations of the machine" and expect it to "reason or think as people do." 46
47 Addressing this challenge involves finding a common language that is simple enough for non-technical stakeholders 47
48 and customizing explanations for different audiences. Another challenge is generating robust and satisfactory ex- 48
49 planations, which can be difficult due to a lack of ground truths, poor model performance, and the black-box nature 49
50 of models. 50
51 51

Task	Example 1			Example 2		
	Representation	Cite	Feedback	Representation	Cite	Feedback
Entity extraction	Words, attention score visualization	TMN [117]	⊕ ⊕ ⊕ ×	Training instances	Instance-based span [119]	⊕ ⊕ × ×
Relation extraction	Words	D-REX [121]	⊕ ⊕ × × ×	Logic rules	LogiRE [127]	✓ ✓ ⊕ ⊕ ×
Entity linking	Attribution scores and their visualization	LEMON [102]	× ×	Entity resolution rules	SystemER [183]	✓ ✓
Link prediction	Reasoning path	[136]	✓ ⊕	Training triples	Kelpie [96]	× ×
Inconsistency detection	Triples and their visualization	Abstraction [181]	× ×		/	

Table 7

User acceptance count of explainable examples, ‘✓’ indicates vote for acceptable explanations, ‘×’ indicates vote for unsatisfactory explanations, and ‘⊕’ indicates vote for explanations that are somewhat reasonable but not fully trustworthy. For inconsistency detection, only one example is provided.

4.3. Gaps and Challenges in Explainable KGC Solutions and Practical Usage

4.3.1. Use Cases from Interview Study

What are practical use cases of XAI models? We first compared the use cases in Section 3.1.2 with those collected from the interview study. We found that the use cases in Section 3.1.2 were largely reflected through the interview study, which also provided new insights and additional use cases. The most prominent use case, highlighted by 10 participants, is understanding the model output and its inner workings. This includes providing supporting evidence, mapping results to the original input, and explaining how the models generate the output. This aligns with the previously identified use case of understanding performance and contributing factors. The second common use case, mentioned by 5 participants, is debugging models and assisting in rectifying and adjusting them. This extends the previous use case of model debugging by indicating where the machine fails or is unstable, identifying systematic error patterns and problematic parts of data sets, and understanding mistakes and errors and their causes.

Two novel use cases were identified from the interviews. The first involves **enhancing human-machine interactions** by facilitating human involvement at various stages of the pipeline and providing effective interaction with the models. Participant A emphasized that having explainability can streamline the workflow, stating that "the more explainable the models are, the less human intervention is required" during model deployment. To this end, clear and informative explanations play a crucial role in bridging the knowledge gap among different stakeholders, ensuring a shared understanding at the right level. Additionally, simplifying the reuse and sharing of results and pipelines among stakeholders and other technical experts is crucial. This is similar to the model update use case in Section 3.1.2, as reusing the pipeline in other processes also involves feeding new data into the pipeline and explaining the differences. By mapping the works discussed in the literature review to this use case, we found that human-in-the-loop approaches, including myDIG [187], SystemER [183], and TuneR [184], align perfectly with this context. Another novel use case that emerged from the interview study is **uncovering new and previously unnoticed insights**. This involves explaining how unexpected (but not necessarily wrong) results are obtained and offering additional details or information that may be overlooked when humans perform the same tasks. Among the works reviewed in 4.1.1, rule-based explanations, including those presented in [108, 129, 183, 184] for entity resolution and [154, 160] for link prediction, demonstrate notable potential for contributing to this use case.

4.3.2. XAI Example Discussion

Do current explainable solutions meet the requirements for practical use cases? During the example discussion session, participants provided feedback on various tasks: 5 commented on relation extraction, 4 on entity extraction, and 2 each on entity resolution, link prediction, and inconsistency detection (Table 7). **Overall, the examples seldom met participants’ requirements, with only 5 out of 28 feedback responses being positive, while nearly half were negative.** Participants highlighted several concerns and issues regarding the practical adoption of the provided explanations. The primary issue, raised in 8 responses, was that the explanations were not sufficiently informative. This could mean several things: the explanations might only cover one or a few aspects of the results, making them insufficient to fully explain the outcomes. Additionally, the correlation or relevance between the explanation and the output was often weak, rendering the explanations inadequate. For instance, using trigger words from the context

1 to explain entity or relation extraction results might show some relevance but still fail to explain why the models 1
2 produced those specific results instead of others. Participants also noted the complexity of the explanations, which 2
3 made them difficult to understand and evaluate. A specific barrier was the use of technical terms. For example, 3
4 explanations represented in logic rules were found useful but too complex for those without a technical background. 4
5 Participant *M* mentioned that numerical thresholds used in rules were perplexing, and participant *C* expressed concerns 5
6 that logic rules could quickly become overly complex, especially for long and intricate contexts. Once an 6
7 error occurred, it was challenging to pinpoint its source, complicating the validation of the explanations. Additionally, 7
8 visualizations or elements used in visualization-based explanations were not always clearly defined, such as 8
9 numbers or color bars, and the representations were not in a standard language familiar to knowledge engineers or 9
10 lay users. A third concern, raised by 6 responses, was the stability and coverage of the explainable models. Participants 10
11 questioned whether these models could consistently provide reliable explanations for all results. This issue, 11
12 known as coverage, refers to how many cases from the entire set of results can be explained. Participants worried 12
13 that there might not always be an explanation, or an explanation of sufficient quality, for all cases of interest. This 13
14 concern was particularly focused on path-based and attention-based explanations. For path-based explanations, participants 14
15 doubted the reliability of reasoning paths for every link prediction result. For attention-based explanations, 15
16 they were concerned about the models' stability and the possibility of incorrect attention being paid to words, thus 16
17 making the explanations less reliable. 17
18

19 4.4. Requirements and Blueprints for Explainable Approaches 19 20

21 4.4.1. Requirements for Explanations 21

22 *What are characteristics of an explainable method that knowledge engineers and researchers expected?* From 22
23 the example discussions, we identified two key requirements. First, four participants emphasized the need for a 23
24 **confidence indicator** for models, explaining how confident the models are when producing results. This indicator 24
25 should truthfully reflect the model's confidence in both the output and the explanations, and models should have 25
26 the option to acknowledge uncertainty rather than provide incorrect answers or hallucinations. Wrong explanations, 26
27 whether for correct or incorrect results, can bias users' impressions of the explanations, thereby eroding user trust 27
28 in the model. Confidence indicators can reduce such cases and facilitate better human-machine collaboration by 28
29 highlighting uncertain instances where human intervention is necessary. As participant *E* noted, when generating 29
30 explanations for relation extraction, there should be "an option of like they don't know the relationships or they 30
31 cannot predict this relationship based on the current context." Similarly, participant *F* stated, "we would like to 31
32 make sure that the machine says it doesn't know when it doesn't know." 32
33

34 Secondly, the representation of explanations largely depends on the task and user. Although explanation formats 34
35 like visualization and logic rules received varying levels of acceptance, the most acceptable representation for participants 35
36 was **natural language**. 11 out of 13 participants preferred natural language explanations, either alone or 36
37 combined with other representations such as visualization and logic rules. The main concern with visualizations 37
38 and logic rules was semantic grounding, meaning the use of clearly defined language that ensures users understand 38
39 the underlying semantics. Participants noted that visual notations without clear definitions are "prone to ambiguity" 39
40 "The same concern was raised for natural language explanations, as ambiguity can create challenges for human 40
41 understanding of the model, thereby complicating human-machine interaction." 41

42 From the requirement elicitation questions, we also identified two common requirements from participants more 42
43 directly. First, 4 participants highlighted that the type of information users most require in explanations is **contextual** 43
44 **information**. They want explainable models capable of tracing which inputs or intermediate outputs led to a certain 44
45 result. As participant *D* mentioned, the ability to "pinpoint" refers to "identifying the minimal relevant information 45
46 that a user needs to understand the result and the problem." This ability to map outputs to inputs can establish 46
47 the basic trustworthiness of explainable models. Participant *I* added, "the sources from where the explanations are 47
48 given are crucial to users." For example, in link prediction tasks, when the input is knowledge graphs, the required 48
49 information includes which triples are used to derive the new one. Generally, this involves providing contextualized 49
50 information with input data and knowledge graphs, reflecting the relevance between input and output, and the 50
51 relevance between specific components or steps and the output. 51

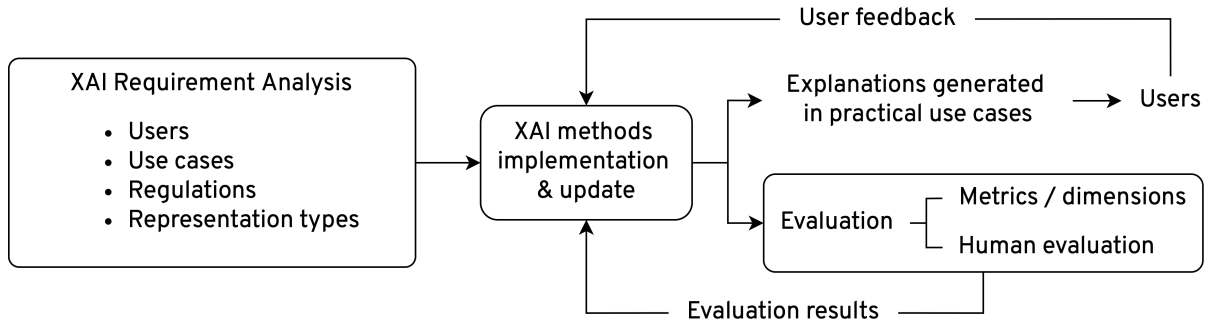


Fig. 6. The XAI methods design blueprint.

Moreover, four participants envisioned a solution involving a "hybrid pipeline" where people and machines work in cooperation, providing ways of **interaction** such as (1) machine-generated explanations for people to understand performance, corner or uncertain cases, etc.; and (2) a feedback loop for people to provide feedback on explanations to the system, explaining why something is right or wrong, or directly giving explanations, so that the machine can learn from this feedback and adjust itself. Follow that line, two participants specifically mentioned the need for iterative explanations in a conversational form. As participant *M* suggested, such "dynamic" explanations could extend over several rounds, allowing users to "keep asking for more depth if they see the need." Similarly, participant *G* believed that the explainable model should be able to select appropriate explanations based on the specific use case and input data, such as rules-based, similarity-based, or visualization-based explanations.

4.4.2. Explanation Design

Based on the results from the literature review and interview study, we propose several guidelines and combine them into blueprints for designing explainable solutions in knowledge engineering tasks that are usable and trustworthy to target users, as shown in Figure 6.

The first step in designing explainable models involves XAI requirement analysis, which collects design insights and creates goals for explanations. Several factors must be carefully considered and investigated to capture the scope and objectives of explainable models.

The most important factor is the **users** who consume the explanations. As participant *A* noted, "designing the system with the users in mind" and "users are the central component." In the context of the knowledge graph lifecycle, these users can be stakeholders, domain experts, knowledge engineers, etc. From the literature review, only 10 works explicitly described the intended users who engage with the generated explanations, as well as their background information. For instance, xERTE [139] reported the background information of respondents involved in the evaluation of explanations, including their education levels. TuneR [184] specifies that the tool is designed to support developers, including coders, scientists, and domain experts. And from the interview study, it is evident that the design of explanations should consider the users' level of understanding and interest in the technical details. Therefore, user analysis should include investigating their background, particularly their technical expertise and domain knowledge, and their expectations for the explanations, including the level of technical detail they require. Furthermore, if the consumers of explanations are involved in collaborative work with machines, understanding how they consume explanations and interact with models is crucial.

The second part of XAI requirement analysis focuses on the **use cases** of explanations. We identified six use cases, each requiring explanations to focus on different aspects. Developers need to decide which practical use cases the explanations will serve, which may extend beyond the six identified. For example, if explanations are used for model debugging and adjustment, they should provide details on inner workings and contributing factors to help identify error sources. A confidence indicator, as mentioned in Section 4.4.1, is also useful. If the goal is to enhance human-AI interactions, it is recommended to design mechanisms for providing explainability at multiple stages through collaboration and a good feedback loop to personalize the model's output based on user input.

The third factor is the **representation** of explanations, which primarily depends on the task and its related input and output (datatype, modality, or property) and user needs. For example, in knowledge extraction tasks where

1 the input modality is text, the context of the original input might be useful (though not necessarily sufficient) as 1
2 supporting evidence to show the relevance of the output. This needs to be combined with user analysis to determine 2
3 what "language" the user speaks, such as description logic, natural language, images, etc. 3

4 Other factors, such as AI regulations mentioned in the Introduction, may also need to be considered in the require- 4
5 ment analysis. Moreover, this list of factors can be expanded based on actual scenarios. Such requirement analysis 5
6 guides the selection and implementation of XAI methods, ensuring they meet the requirements of practical use 6
7 cases. After implementing XAI methods, the workflow involves iterative loops for maintaining and continuously 7
8 improving the methods. One loop focuses on the evaluation and assessment of explanations [202]. Evaluation meth- 8
9 ods should go beyond anecdotal evidence, selecting appropriate metrics or designing evaluation paradigms. Another 9
10 iterative loop, derived from the "hybrid pipeline" requirements in Section 4.4.1, aims to improve explainable mod- 10
11 els and explanations in practical scenarios. Users who consume the explanations provide feedback and example 11
12 explanations, which can be used in various ways to enhance the XAI model. This includes creating datasets of ex- 12
13 planations for training and fine-tuning XAI models, providing few-shot examples, or even abstracting improvement 13
14 directions for architecture-level adjustments. 14

15 16 17 **5. Conclusion and Future Work** 17

18 19 *5.1. Conclusion* 19

20 21 In this paper, we adopted a mixed methodology, conducting a literature review on explainable methods within the 21
22 domain of knowledge graph construction and an interview study on the same topic with 13 participants to capture 22
23 how XAI methods support knowledge engineering. We performed the analysis in three dimensions, tasks related 23
24 to knowledge graph construction, the taxonomy of XAI methods, and the use cases of XAI methods in knowledge 24
25 graph construction. We observed that the most effort has been directed towards automation and explainability in 25
26 entity extraction, relation extraction, entity linking, and link prediction. Additionally, we considered the use cases 26
27 in explainable automatic knowledge graph construction, such as ML model selecting and building, ML model de- 27
28 bugging, understanding performance and contributing factors, and managing updates. The interview study largely 28
29 corroborated the considered use cases, adding new insights and highlighting additional use cases, including enhanc- 29
30 ing human-machine interactions and providing new insights from unexpected results. We found that the reviewed 30
31 models primarily focused on explaining the performance and contributing factors to the outcome while neglecting 31
32 other use cases, such as error detection and correction, which could help establish trust with users. The interview 32
33 study revealed that while current knowledge engineering tools exhibit varying degrees of automation and under- 33
34 standing, significant challenges remain in data provenance, evaluation methods, and providing clear explanations to 34
35 stakeholders. The current explainable solutions often fell short of participants' requirements, with concerns about 35
36 their informativeness, complexity, and reliability. These insights established a foundational understanding of criti- 36
37 cal transparency factors, enabling the development of a blueprint for designing explainable methods for knowledge 37
38 engineering tasks. 38

39 40 *5.2. Future Work* 40

41 42 We identified five future directions for research on explainable automatic knowledge graph construction. First, 42
43 going back to prior literature on knowledge engineering methodologies [59, 60, 203, 204], there are many **tasks** 43
44 **and activities where automation remains an exception**. Aside from the tasks in Figure 3, there is an opportunity 44
45 to think about other ways for AI assistance to add value: for instance, one design principle of KGs is that they are 45
46 meant to integrate across multiple sources and be able to tackle evolving requirements. Reusing existing schemas 46
47 or ontologies can help with interoperability, but the task of finding or assessing an ontology for reuse is still mostly 47
48 manual. At the other end of the lifecycle, documenting KGs can help with maintenance and reuse, and advances 48
49 in generative AI make it a chief candidate for automation. While we found a range of explainable link prediction 49
50 approaches, it would be useful to dive deeper into this sub-field to understand the extent to which these different 50
51 approaches solve common concerns around the quality of KGs. One difference between representing knowledge in 51

a KG and a machine-learning model is that a KG can provide guarantees about the validity of the information, its provenance, its currency, etc. upon retrieval. However, this is predicated by KGs being regularly audited according to these and other quality dimensions and improved. Link prediction is one way to do this, alongside many others, e.g., debiasing [16]. Furthermore, while knowledge acquisition is generally well represented in the literature, a lot of work focuses on text rather than other data modalities, which is a concern in many KG application areas, e.g., enterprise data management (which needs to work with structured data) or cultural heritage (where a lot of domain data is neither text nor numbers).

Second, as we noted earlier, the fewest of approaches look at **the human-in-the-loop aspects of KG construction**, including human agency and oversight, feedback, etc [191] and **the integration of the developed models into established knowledge-engineering practices**. While there is a lot of work in human-AI interaction and interactive ML in the HCI community, they tend to focus so far on simpler ML models and different applications that the knowledge production scenarios we are interested in. One exception is the work on ORES [205], a participatory ML system used in Wikipedia and Wikidata (a large open-source KG). However, the Wikidata KG construction process is unique because it is community-based, with more than 24K active contributors¹⁵ who receive AI assistance for distinct tasks such as vandalism detection and consistency checks. We need to follow their example to develop the same types of workflows and tools for other KG construction scenarios - in most cases, these involve much smaller teams and different tool environments. The majority of existing integrated development environments (IDEs) for KGs (e.g., PoolParty¹⁶, data.world¹⁷, Protégé¹⁸) assume KGs are mostly built manually, with some basic automation to speed-up routine tasks like translating node labels or creating documentation from node and edge descriptions. LLMs like ChatGPT offer chances to develop novel KG editing tools and interactions, allowing people to interact with their AI assistants via natural language and ensuring transparency. Meanwhile, developers working with KGs require KG-related process blueprints that utilize AI algorithms and adhere to AI regulations for creating downstream applications.

Thirdly, our research flagged the need for **better evaluations on explanations**, which encompasses metrics, benchmarks, and datasets, as well as toolkits and guidance for conducting studies that assess how effective the explanations supplied in KG construction tasks are as proxies and enablers for transparent and hence trusted KGs.

Fourthly, our research revealed an imbalance in the distribution of use cases identified in the study. There was a strong emphasis on understanding the inner workings, performance, and contributing factors of models, while relatively few efforts were made to address other use cases also demanded by the community, such as model debugging, model updating, and human-AI interaction. However, our example discussions indicated that the reviewed explanations often failed to meet these requirements, and participants expressed low confidence in using them in their work or providing them to users. A future direction, reflected in our study and requested by the community, involves **adapting current explainable methods to representations and formats that are reusable across multiple use cases**.

Finally, although our research provided a blueprint for designing XAI methods, **practical applications and verification of the blueprint** are missing. Given the different use cases and groups of stakeholders in a knowledge engineering project, several details can be enriched. For instance, parts of the blueprint, such as the user feedback loop, can be refined. Future work could investigate what formats, workflows, and feedback frequencies best prompt users to provide high-quality explanations efficiently.

References

- [1] A. Hogan, E. Blomqvist, M. Cochez, C. d'Amato, G. de Melo, C. Gutierrez, S. Kirrane, J.E.L. Gayo, R. Navigli, S. Neumaier et al., *Knowledge Graphs*, Synthesis Lectures on Data, Semantics, and Knowledge, Morgan & Claypool Publishers, 2021. ISBN 9781636392363. <https://books.google.co.uk/books?id=hJ1NEAAAQBAJ>.

¹⁵<https://www.wikidata.org/wiki/Wikidata:Statistics>

¹⁶<https://www.poolparty.biz/>

¹⁷<https://data.world/>

¹⁸<https://protege.stanford.edu/>

- [2] J. Chen, Y. Geng, Z. Chen, J.Z. Pan, Y. He, W. Zhang, I. Horrocks and H. Chen, Zero-shot and Few-shot Learning with Knowledge Graphs: A Comprehensive Survey, *CoRR abs/2112.10006* (2021). <https://arxiv.org/abs/2112.10006>.
- [3] P.S.H. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W. Yih, T. Rocktäschel, S. Riedel and D. Kiela, Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks, *CoRR abs/2005.11401* (2020). <https://arxiv.org/abs/2005.11401>.
- [4] L. Yang, H. Chen, Z. Li, X. Ding and X. Wu, Give Us the Facts: Enhancing Large Language Models with Knowledge Graphs for Fact-aware Language Modeling, 2024.
- [5] I. Tiddi and S. Schlobach, Knowledge graphs as tools for explainable machine learning: A survey, *Artificial Intelligence* **302** (2022), 103627. doi:<https://doi.org/10.1016/j.artint.2021.103627>. <https://www.sciencedirect.com/science/article/pii/S0004370221001788>.
- [6] J. Sequeda and O. Lassila, *Designing and Building Enterprise Knowledge Graphs*, Springer International Publishing, 2021. ISSN 2691-2031. ISBN 9783031019166. doi:10.1007/978-3-031-01916-6.
- [7] T.B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D.M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever and D. Amodei, Language Models are Few-Shot Learners, *CoRR abs/2005.14165* (2020). <https://arxiv.org/abs/2005.14165>.
- [8] OpenAI, GPT-4 Technical Report, *arXiv preprint arXiv:2303.08774* (2023).
- [9] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, A. Rodriguez, A. Joulin, E. Grave and G. Lample, LLaMA: Open and Efficient Foundation Language Models, 2023.
- [10] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, D. Bikel, L. Blecher, C.C. Ferrer, M. Chen, G. Cucurull, D. Esiobu, J. Fernandez, J. Fu, W. Fu, B. Fuller, C. Gao, V. Goswami, N. Goyal, A. Hartshorn, S. Hosseini, R. Hou, H. Inan, M. Kardas, V. Kerkez, M. Khabsa, I. Kloumann, A. Korenev, P.S. Koura, M.-A. Lachaux, T. Lavril, J. Lee, D. Liskovich, Y. Lu, Y. Mao, X. Martinet, T. Mihaylov, P. Mishra, I. Molybog, Y. Nie, A. Poulton, J. Reizenstein, R. Rungta, K. Saladi, A. Schelten, R. Silva, E.M. Smith, R. Subramanian, X.E. Tan, B. Tang, R. Taylor, A. Williams, J.X. Kuan, P. Xu, Z. Yan, I. Zarov, Y. Zhang, A. Fan, M. Kambadur, S. Narang, A. Rodriguez, R. Stojnic, S. Edunov and T. Scialom, Llama 2: Open Foundation and Fine-Tuned Chat Models, 2023.
- [11] J.Z. Pan, S. Razniewski, J.-C. Kalo, S. Singhanian, J. Chen, S. Dietze, H. Jabeen, J. Omeljanenko, W. Zhang, M. Lissandrini, R. Biswas, G. de Melo, A. Bonifati, E. Vakaj, M. Dragoni and D. Graux, Large Language Models and Knowledge Graphs: Opportunities and Challenges, *Transactions on Graph Data and Knowledge* **1**(1) (2023), 2:1–2:38. doi:10.4230/TGDK.1.1.2.
- [12] G.e. Tamašauskaitė and P. Groth, Defining a Knowledge Graph Development Process Through a Systematic Review, *ACM Trans. Softw. Eng. Methodol.* **32**(1) (2023). doi:10.1145/3522586.
- [13] G. Weikum, L. Dong, S. Razniewski and F.M. Suchanek, Machine Knowledge: Creation and Curation of Comprehensive Knowledge Bases, *CoRR abs/2009.11564* (2020). <https://arxiv.org/abs/2009.11564>.
- [14] C.T. Wolf, From Knowledge Graphs to Knowledge Practices: On the Need for Transparency and Explainability in Enterprise Knowledge Graph Applications, in: *Knowledge Graph Bias Workshop*, 2020.
- [15] D. Abián, A. Meroño-Peñuela and E. Simperl, An Analysis of Content Gaps Versus User Needs in the Wikidata Knowledge Graph, in: *The Semantic Web—ISWC 2022: 21st International Semantic Web Conference, Virtual Event, October 23–27, 2022, Proceedings*, Springer, 2022, pp. 354–374.
- [16] J. Fisher, A. Mittal, D. Palfrey and C. Christodoulopoulos, Debiasing knowledge graph embeddings, in: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020, pp. 7332–7345.
- [17] L. Beyer, O.J. Hénaff, A. Kolesnikov, X. Zhai and A.v.d. Oord, Are we done with imagenet?, *arXiv preprint arXiv:2006.07159* (2020).
- [18] D. Danks and A.J. London, Algorithmic Bias in Autonomous Systems., in: *Ijcai*, Vol. 17, 2017, pp. 4691–4697.
- [19] S. Hooker, Moving beyond “algorithmic bias is a data problem”, *Patterns* **2**(4) (2021), 100241.
- [20] T. Panch, H. Mattie and R. Atun, Artificial intelligence and algorithmic bias: implications for health systems, *Journal of global health* **9**(2) (2019).
- [21] M.T. Ribeiro, T. Wu, C. Guestrin and S. Singh, Beyond accuracy: Behavioral testing of NLP models with CheckList, *arXiv preprint arXiv:2005.04118* (2020).
- [22] B. Zhang, A. Meroño Peñuela and E. Simperl, Towards Explainable Automatic Knowledge Graph Construction with Human-in-the-Loop, in: *HAI 2023: Augmenting Human Intellect*, IOS Press, 2023, pp. 274–289. doi:10.3233/FAIA230091.
- [23] P. Groth, E. Simperl, M. van Erp and D. Vrandečić, Knowledge Graphs and their Role in the Knowledge Engineering of the 21st Century (Dagstuhl Seminar 22372), *Dagstuhl Reports* **12**(9) (2023), 60–120. doi:10.4230/DagRep.12.9.60. <https://drops.dagstuhl.de/opus/volltexte/2023/17810>.
- [24] J.D. Lee and K.A. See, Trust in Automation: Designing for Appropriate Reliance, *Human Factors* **46**(1) (2004), 50–80, PMID: 15151155. doi:10.1518/hfes.46.1.50_30392. https://doi.org/10.1518/hfes.46.1.50_30392.
- [25] D. Vrandečić and M. Krötzsch, Wikidata: A Free Collaborative Knowledgebase, *Commun. ACM* **57**(10) (2014), 78–85. doi:10.1145/2629489.
- [26] J. Lehmann, R. Isele, M. Jakob, A. Jentzsch, D. Kontokostas, P.N. Mendes, S. Hellmann, M. Morsey, P. Van Kleef, S. Auer et al., DBpedia - A large-scale, multilingual knowledge base extracted from Wikipedia, *Semantic Web* **6**(2) (2015), 167–195.
- [27] T. Pellissier Tanon, G. Weikum and F. Suchanek, YAGO 4: A Reason-able Knowledge Base, in: *The Semantic Web*, A. Harth, S. Kirrane, A.-C. Ngonga Ngomo, H. Paulheim, A. Rula, A.L. Gentile, P. Haase and M. Cochez, eds, Springer International Publishing, Cham, 2020, pp. 583–596. ISBN 978-3-030-49461-2.
- [28] R. Speer, J. Chin and C. Havasi, ConceptNet 5.5: An Open Multilingual Graph of General Knowledge, *CoRR abs/1612.03975* (2016). <http://arxiv.org/abs/1612.03975>.

- [29] M. van Bakkum, M. de Boer, F. van Harmelen, A. Meyer-Vitali and A.t. Teije, Modular design patterns for hybrid learning and reasoning systems: a taxonomy, patterns and use cases, *Applied Intelligence* **51**(9) (2021), 6528–6546.
- [30] A. Breit, L. Waltersdorfer, F.J. Ekaputra, M. Sabou, A. Ekelhart, A. Iana, H. Paulheim, J. Portisch, A. Revenko, A.t. Teije and F. van Harmelen, Combining Machine Learning and Semantic Web: A Systematic Mapping Study, *ACM Comput. Surv.* (2023), Just Accepted. doi:10.1145/3586163.
- [31] F. Poursabzi-Sangdeh, D.G. Goldstein, J.M. Hofman, J.W. Vaughan and H.M. Wallach, Manipulating and Measuring Model Interpretability, *CoRR abs/1802.07810* (2018). <http://arxiv.org/abs/1802.07810>.
- [32] A. Smith-Renner, R. Fan, M. Birchfield, T. Wu, J. Boyd-Graber, D.S. Weld and L. Findlater, No Explainability without Accountability: An Empirical Study of Explanations and Feedback in Interactive ML, in: *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, CHI '20, Association for Computing Machinery, New York, NY, USA, 2020, pp. 1–13–. ISBN 9781450367080. doi:10.1145/3313831.3376624.
- [33] X. Wang and M. Yin, Effects of Explanations in AI-Assisted Decision Making: Principles and Comparisons **12**(4) (2022). doi:10.1145/3519266.
- [34] S. Amershi, D. Weld, M. Vorvoreanu, A. Fourney, B. Nushi, P. Collisson, J. Suh, S. Iqbal, P.N. Bennett, K. Inkpen et al., Guidelines for human-AI interaction, in: *Proceedings of the 2019 chi conference on human factors in computing systems*, 2019, pp. 1–13.
- [35] J.A. Fails and D.R. Olsen Jr, Interactive machine learning, in: *Proceedings of the 8th international conference on Intelligent user interfaces*, 2003, pp. 39–45.
- [36] C. Sarasua, E. Simperl, N.F. Noy, A. Bernstein and J.M. Leimeister, Crowdsourcing and the semantic web: A research manifesto, *Human Computation* **2**(1) (2015).
- [37] E. Simperl, R. Cuel and M. Stein, Incentive-centric semantic web application engineering, *Synthesis Lectures on the Semantic Web: Theory and Technology* **3**(1) (2013), 1–117.
- [38] U. Ehsan, Q.V. Liao, M. Muller, M.O. Riedl and J.D. Weisz, Expanding explainability: Towards social transparency in AI systems, in: *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, 2021, pp. 1–19.
- [39] D. Kaur, S. Uslu, K.J. Rittichier and A. Durresi, Trustworthy Artificial Intelligence: A Review, *ACM Comput. Surv.* **55**(2) (2022). doi:10.1145/3491209.
- [40] S. Larsson and F. Heintz, Transparency in artificial intelligence, *Internet Policy Review* **9**(2) (2020).
- [41] G. Schwalbe and B. Finzel, XAI Method Properties: A (Meta-)study, *CoRR abs/2105.07190* (2021). <https://arxiv.org/abs/2105.07190>.
- [42] M.T. Ribeiro, S. Singh and C. Guestrin, ‘Why Should I Trust You?’: Explaining the Predictions of Any Classifier, *CoRR abs/1602.04938* (2016). <http://arxiv.org/abs/1602.04938>.
- [43] S.M. Lundberg and S.-I. Lee, A Unified Approach to Interpreting Model Predictions, in: *Advances in Neural Information Processing Systems 30*, I. Guyon, U.V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan and R. Garnett, eds, Curran Associates, Inc., 2017, pp. 4765–4774. <http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf>.
- [44] P. Hase and M. Bansal, Evaluating Explainable AI: Which Algorithmic Explanations Help Users Predict Model Behavior?, in: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, Association for Computational Linguistics, Online, 2020, pp. 5540–5552. doi:10.18653/v1/2020.acl-main.491. <https://aclanthology.org/2020.acl-main.491>.
- [45] D. Minh, H.X. Wang, Y.F. Li and T.N. Nguyen, Explainable Artificial Intelligence: A Comprehensive Review, *Artif. Intell. Rev.* **55**(5) (2022), 3503–3568–. doi:10.1007/s10462-021-10088-y.
- [46] S. Mohseni, N. Zarei and E.D. Ragan, A Multidisciplinary Survey and Framework for Design and Evaluation of Explainable AI Systems, *ACM Trans. Interact. Intell. Syst.* **11**(3–4) (2021). doi:10.1145/3387166.
- [47] G. Vilone and L. Longo, Explainable Artificial Intelligence: a Systematic Review, *CoRR abs/2006.00093* (2020). <https://arxiv.org/abs/2006.00093>.
- [48] A.B. Arrieta, N.D. Rodríguez, J.D. Ser, A. Bannetot, S. Tabik, A. Barbado, S. García, S. Gil-Lopez, D. Molina, R. Benjamins, R. Chatila and F. Herrera, Explainable Artificial Intelligence (XAI): Concepts, Taxonomies, Opportunities and Challenges toward Responsible AI, *CoRR abs/1910.10045* (2019). <http://arxiv.org/abs/1910.10045>.
- [49] M. Danilevsky, K. Qian, R. Aharonov, Y. Katsis, B. Kawas and P. Sen, A Survey of the State of Explainable AI for Natural Language Processing, *CoRR abs/2010.00711* (2020). <https://arxiv.org/abs/2010.00711>.
- [50] T. Miller, Explanation in artificial intelligence: Insights from the social sciences, *Artificial Intelligence* **267** (2019), 1–38. doi:<https://doi.org/10.1016/j.artint.2018.07.007>. <https://www.sciencedirect.com/science/article/pii/S0004370218305988>.
- [51] Y. Rong, T. Leemann, T.-t. Nguyen, L. Fiedler, P. Qian, V. Unhelkar, T. Seidel, G. Kasneci and E. Kasneci, Towards Human-centered Explainable AI: User Studies for Model Explanations, 2022.
- [52] A.D. Preece, D. Harborne, D. Braines, R. Tomsett and S. Chakraborty, Stakeholders in Explainable AI, *CoRR abs/1810.00184* (2018). <http://arxiv.org/abs/1810.00184>.
- [53] G. Ras, M. van Gerven and P. Haselager, *Explanation Methods in Deep Learning: Users, Values, Concerns and Challenges*, in: *Explainable and Interpretable Models in Computer Vision and Machine Learning*, H.J. Escalante, S. Escalera, I. Guyon, X. Baró, Y. Güçlütürk, U. Güçlü and M. van Gerven, eds, Springer International Publishing, Cham, 2018, pp. 19–36. ISBN 978-3-319-98131-4. doi:10.1007/978-3-319-98131-4_2.
- [54] Q.V. Liao, D. Gruen and S. Miller, Questioning the AI: Informing Design Practices for Explainable AI User Experiences, in: *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, CHI '20, Association for Computing Machinery, New York, NY, USA, 2020, pp. 1–15–. ISBN 9781450367080. doi:10.1145/3313831.3376590.

- [55] S. Dhanorkar, C.T. Wolf, K. Qian, A. Xu, L. Popa and Y. Li, Who Needs to Know What, When?: Broadening the Explainable AI (XAI) Design Space by Looking at Explanations Across the AI Lifecycle, in: *Designing Interactive Systems Conference 2021*, DIS '21, Association for Computing Machinery, New York, NY, USA, 2021, pp. 1591–1602-. ISBN 9781450384766. doi:10.1145/3461778.3462131.
- [56] S.S.Y. Kim, E.A. Watkins, O. Russakovsky, R. Fong and A. Monroy-Hernández, "Help Me Help the AI": Understanding How Explainability Can Support Human-AI Interaction, in: *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, CHI '23, Association for Computing Machinery, New York, NY, USA, 2023. ISBN 9781450394215. doi:10.1145/3544548.3581001.
- [57] M. Nauta, J. Trienes, S. Pathak, E. Nguyen, M. Peters, Y. Schmitt, J. Schlötterer, M. van Keulen and C. Seifert, From Anecdotal Evidence to Quantitative Evaluation Methods: A Systematic Review on Evaluating Explainable AI, *CoRR* **abs/2201.08164** (2022). <https://arxiv.org/abs/2201.08164>.
- [58] M. Chromik and M. Schuessler, A Taxonomy for Human Subject Evaluation of Black-Box Explanations in XAI, *ExSS-ATEC@IUI 1* (2020).
- [59] G. Schreiber, *Knowledge Engineering and Management: The CommonKADS Methodology*, A Bradford book, MIT Press, 2000. ISBN 9780262193009. https://books.google.co.uk/books?id=HIXOW_1fsIEC.
- [60] R. Studer, V.R. Benjamins and D. Fensel, Knowledge engineering: Principles and methods, *Data & knowledge engineering* **25**(1–2) (1998), 161–197.
- [61] D. Fensel, U. Simsek, K. Angele, E. Huaman, E. Kärle, O. Panasiuk, I. Toma, J. Umbrich and A. Wahler, *Knowledge graphs*, Springer, 2020.
- [62] M.K. Sarker, L. Zhou, A. Eberhart and P. Hitzler, Neuro-symbolic artificial intelligence, *AI Communications* **34**(3) (2021), 197–209.
- [63] P. Schneider, T. Schopf, J. Vladika, M. Galkin, E. Simperl and F. Matthes, A Decade of Knowledge Graphs in Natural Language Processing: A Survey, in: *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, Association for Computational Linguistics, Online only, 2022, pp. 601–614. <https://aclanthology.org/2022.aacl-main.46>.
- [64] H. Ye, N. Zhang, H. Chen and H. Chen, Generative Knowledge Graph Construction: A Review, *CoRR* **abs/2210.12714** (2022). doi:10.48550/arXiv.2210.12714.
- [65] B. Zhang, V.A. Carriero, K. Schreiberhuber, S. Tsaneva, L.S. González, J. Kim and J. de Berardinis, OntoChat: a Framework for Conversational Ontology Engineering using Language Models, *arXiv preprint arXiv:2403.05921* (2024).
- [66] E. Simperl and M. Luczak-Rösch, Collaborative ontology engineering: a survey, *The Knowledge Engineering Review* **29**(1) (2014), 101–131-. doi:10.1017/S0269888913000192.
- [67] U. Simsek, E. Kärle, K. Angele, E. Huaman, J. Opendplatz, D. Sommer, J. Umbrich and D. Fensel, A Knowledge Graph Perspective on Knowledge Engineering, *SN Comput. Sci.* **4**(1) (2022). doi:10.1007/s42979-022-01429-x.
- [68] F. van Harmelen and A. ten Teije, A Boxology of Design Patterns for Hybrid Learning and Reasoning Systems, *Journal of Web Engineering* **18**(1–3) (2019). <https://journals.riverpublishers.com/index.php/JWE/article/view/3175>.
- [69] H.F. Witschel, C. Pande, A. Martin, E. Laurenzi and K. Hinkelmann, *Visualization of Patterns for Hybrid Learning and Reasoning with Human Involvement*, in: *New Trends in Business Information Systems and Technology: Digital Innovation and Digital Business Transformation*, R. Dornberger, ed., Springer International Publishing, Cham, 2021, pp. 193–204. ISBN 978-3-030-48332-6. doi:10.1007/978-3-030-48332-6_13. https://doi.org/10.1007/978-3-030-48332-6_{_}13.
- [70] C.W. Holsapple and K.D. Joshi, A Collaborative Approach to Ontology Design, *Commun. ACM* **45**(2) (2002), 42–47-. doi:10.1145/503124.503147.
- [71] S. Auer and H. Herre, RapidOWL — An Agile Knowledge Engineering Methodology, in: *Perspectives of Systems Informatics*, I. Virbitskaite and A. Voronkov, eds, Springer Berlin Heidelberg, Berlin, Heidelberg, 2007, pp. 424–430. ISBN 978-3-540-70881-0.
- [72] S. Braun, A.P. Schmidt, A. Walter, G. Nagypál and V. Zacharias, Ontology Maturing: a Collaborative Web 2.0 Approach to Ontology Engineering, in: *CKC*, 2007.
- [73] C. Debruyne, T.-K. Tran and R. Meersman, Grounding Ontologies with Social Processes and Natural Language, *Journal on Data Semantics* **2**(2–3) (2013), 89–118. doi:10.1007/s13740-013-0023-3.
- [74] K. Kotis and A. Vouros, Human-Centered Ontology Engineering: The HCOME Methodology, *Knowl. Inf. Syst.* **10**(1) (2006), 109–131-. doi:10.1007/s10115-005-0227-4.
- [75] D. Vrandečić, H.S. Pinto, C. Tempich and Y. Sure-Vetter, The DILIGENT knowledge processes, *Journal of Knowledge Management* **9** (2005), 85–96.
- [76] A. de Moor, P. De Leenheer and R. Meersman, DOGMA-MESS: A Meaning Evolution Support System for Interorganizational Ontology Engineering, in: *Conceptual Structures: Inspiration and Application*, H. Schärfe, P. Hitzler and P. Øhrstrøm, eds, Springer Berlin Heidelberg, Berlin, Heidelberg, 2006, pp. 189–202. ISBN 978-3-540-35902-9.
- [77] N. Guarino and C.A. Welty, *An Overview of OntoClean*, in: *Handbook on Ontologies*, S. Staab and R. Studer, eds, Springer Berlin Heidelberg, Berlin, Heidelberg, 2004, pp. 151–171. ISBN 978-3-540-24750-0. doi:10.1007/978-3-540-24750-0_8.
- [78] M. Poveda-Villalón, A. Gómez-Pérez and M.C. Suárez-Figueroa, OOPS! (Ontology Pitfall Scanner!): An On-Line Tool for Ontology Evaluation, *Int. J. Semant. Web Inf. Syst.* **10**(2) (2014), 7–34-. doi:10.4018/ijswis.2014040102.
- [79] E.F. Kendall and D.L. McGuinness, *Requirements and Use Cases*, in: *Ontology Engineering*, Springer International Publishing, Cham, 2019, pp. 25–44. ISBN 978-3-031-79486-5. doi:10.1007/978-3-031-79486-5_3. https://doi.org/10.1007/978-3-031-79486-5_{_}3.
- [80] V. Yadav and S. Bethard, A Survey on Recent Advances in Named Entity Recognition from Deep Learning models, *CoRR* **abs/1910.11470** (2019). <http://arxiv.org/abs/1910.11470>.

- [81] Y. Lin, S. Shen, Z. Liu, H. Luan and M. Sun, Neural Relation Extraction with Selective Attention over Instances, in: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Association for Computational Linguistics, Berlin, Germany, 2016, pp. 2124–2133. doi:10.18653/v1/P16-1200. <https://aclanthology.org/P16-1200>.
- [82] Ö. Sevgili, A. Shelmanov, M.Y. Arkhipov, A. Panchenko and C. Biemann, Neural Entity Linking: A Survey of Models based on Deep Learning, *CoRR abs/2006.00575* (2020). <https://arxiv.org/abs/2006.00575>.
- [83] M. Acosta, A. Zaveri, E. Simperl, D. Kontokostas, S. Auer and J. Lehmann, Crowdsourcing Linked Data Quality Assessment, in: *The Semantic Web – ISWC 2013*, H. Alani, L. Kagal, A. Fokoue, P. Groth, C. Biemann, J.X. Parreira, L. Aroyo, N. Noy, C. Welty and K. Janowicz, eds, Springer Berlin Heidelberg, Berlin, Heidelberg, 2013, pp. 260–276. ISBN 978-3-642-41338-4.
- [84] A. Revenko, M. Sabou, A. Ahmeti and M. Schauer, Crowd-Sourced Knowledge Graph Extension: A Belief Revision Based Approach., in: *HCOMP (WIP&Demo)*, 2018.
- [85] Z. Kou, Y. Zhang, D. Zhang and D. Wang, CrowdGraph: A Crowdsourcing Multi-modal Knowledge Graph Approach to Explainable Fauxtography Detection, *Proc. ACM Hum.-Comput. Interact.* **6**(CSCW2) (2022). doi:10.1145/3555178.
- [86] A. Piscopo and E. Simperl, What we talk about when we talk about wikidata quality: a literature survey, in: *Proceedings of the 15th International Symposium on Open Collaboration*, OpenSym '19, Association for Computing Machinery, New York, NY, USA, 2019. ISBN 9781450363198. doi:10.1145/3306446.3340822.
- [87] K. Shenoy, F. Ilievski, D. Garijo, D. Schwabe and P.A. Szekely, A Study of the Quality of Wikidata, *CoRR abs/2107.00156* (2021). <https://arxiv.org/abs/2107.00156>.
- [88] E. Koutsiana, T. Yadav, N. Jain, A. Meroño-Peñuela and E. Simperl, Agreeing and Disagreeing in Collaborative Knowledge Graph Construction: An Analysis of Wikidata, *CoRR abs/2306.11766* (2023). doi:10.48550/ARXIV.2306.11766. <https://doi.org/10.48550/arXiv.2306.11766>.
- [89] R. Qarout, A. Checco, G. Demartini and K. Bontcheva, Platform-related factors in repeatability and reproducibility of crowdsourcing tasks, in: *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, Vol. 7, 2019, pp. 135–143.
- [90] A. Rossi, D. Firmani, A. Matinata, P. Merialdo and D. Barbosa, Knowledge Graph Embedding for Link Prediction: A Comparative Analysis, *CoRR abs/2002.00819* (2020). <https://arxiv.org/abs/2002.00819>.
- [91] P. Cimiano and H. Paulheim, Knowledge graph refinement: A survey of approaches and evaluation methods, *Semant. Web* **8**(3) (2017), 489–508–. doi:10.3233/SW-160218.
- [92] M. Wiegmann, M. Völske, B. Stein and M. Potthast, Language Models as Context-sensitive Word Search Engines, in: *Proceedings of the First Workshop on Intelligent and Interactive Writing Assistants (In2Writing 2022)*, Association for Computational Linguistics, Dublin, Ireland, 2022, pp. 39–45. doi:10.18653/v1/2022.in2writing-1.5. <https://aclanthology.org/2022.in2writing-1.5>.
- [93] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H.W. Chung, C. Sutton, S. Gehrmann, P. Schuh, K. Shi, S. Tsvyashchenko, J. Maynez, A. Rao, P. Barnes, Y. Tay, N. Shazeer, V. Prabhakaran, E. Reif, N. Du, B. Hutchinson, R. Pope, J. Bradbury, J. Austin, M. Isard, G. Gur-Ari, P. Yin, T. Duke, A. Levskaya, S. Ghemawat, S. Dev, H. Michalewski, X. Garcia, V. Misra, K. Robinson, L. Fedus, D. Zhou, D. Ippolito, D. Luan, H. Lim, B. Zoph, A. Spiridonov, R. Sepassi, D. Dohan, S. Agrawal, M. Omernick, A.M. Dai, T.S. Pillai, M. Pellat, A. Lewkowycz, E. Moreira, R. Child, O. Polozov, K. Lee, Z. Zhou, X. Wang, B. Saeta, M. Diaz, O. Firat, M. Catasta, J. Wei, K. Meier-Hellstern, D. Eck, J. Dean, S. Petrov and N. Fiedel, PaLM: Scaling Language Modeling with Pathways, arXiv, 2022. doi:10.48550/ARXIV.2204.02311. <https://arxiv.org/abs/2204.02311>.
- [94] J. Guo, J. Li, D. Li, A.M.H. Tiong, B. Li, D. Tao and S.C.H. Hoi, From Images to Textual Prompts: Zero-shot VQA with Frozen Large Language Models, arXiv, 2022. doi:10.48550/ARXIV.2212.10846. <https://arxiv.org/abs/2212.10846>.
- [95] Y. Gao, Y. Xiong, X. Gao, K. Jia, J. Pan, Y. Bi, Y. Dai, J. Sun, Q. Guo, M. Wang and H. Wang, Retrieval-Augmented Generation for Large Language Models: A Survey, *CoRR abs/2312.10997* (2023). doi:10.48550/ARXIV.2312.10997. <https://doi.org/10.48550/arXiv.2312.10997>.
- [96] A. Rossi, D. Firmani, P. Merialdo and T. Teofili, Explaining Link Prediction Systems Based on Knowledge Graph Embeddings, in: *Proceedings of the 2022 International Conference on Management of Data*, SIGMOD '22, Association for Computing Machinery, New York, NY, USA, 2022, pp. 2062–2075–. ISBN 9781450392495. doi:10.1145/3514221.3517887.
- [97] F. Bianchi, G. Rossiello, L. Costabello, M. Palmonari and P. Minervini, Knowledge Graph Embeddings and Explainable AI, *CoRR abs/2004.14843* (2020). <https://arxiv.org/abs/2004.14843>.
- [98] M.J. Page, D. Moher, P.M. Bossuyt, I. Boutron, T.C. Hoffmann, C.D. Mulrow, L. Shamseer, J.M. Tetzlaff, E.A. Akl, S.E. Brennan, R. Chou, J. Glanville, J.M. Grimshaw, A. Hróbjartsson, M.M. Lalu, T. Li, E.W. Loder, E. Mayo-Wilson, S. McDonald, L.A. McGuinness, L.A. Stewart, J. Thomas, A.C. Tricco, V.A. Welch, P. Whiting and J.E. McKenzie, PRISMA 2020 explanation and elaboration: updated guidance and exemplars for reporting systematic reviews, *BMJ* **372** (2021). doi:10.1136/bmj.n160. <https://www.bmj.com/content/372/bmj.n160>.
- [99] K. Amarasinghe, K.T. Rodolfa, H. Lamba and R. Ghani, Explainable machine learning for public policy: Use cases, gaps, and research directions, *Data & Policy* **5** (2023), e5. doi:10.1017/dap.2023.2.
- [100] S. Chari, O. Seneviratne, M. Ghalwash, S. Shirai, D.M. Gruen, P. Meyer, P. Chakraborty and D.L. McGuinness, Explanation Ontology: A general-purpose, semantic representation for supporting user-centered explanations, *Semantic Web* (2023), 1–31.
- [101] D. Firmani, L. Tanca and R. Torlone, Ethical Dimensions for Data Quality, *J. Data and Information Quality* **12**(1) (2019). doi:10.1145/3362121.
- [102] N. Barlaug, LEMON: Explainable Entity Matching, *CoRR abs/2110.00516* (2021). <https://arxiv.org/abs/2110.00516>.
- [103] J. Wang and Y. Li, Minun: Evaluating Counterfactual Explanations for Entity Matching, in: *Proceedings of the Sixth Workshop on Data Management for End-To-End Machine Learning*, DEEM '22, Association for Computing Machinery, New York, NY, USA, 2022. ISBN 9781450393751. doi:10.1145/3533028.3533304.

- [104] A. Baraldi, F. Del Buono, M. Paganelli and F. Guerra, Landmark Explanation: An Explainer for Entity Matching Models, in: *Proceedings of the 30th ACM International Conference on Information and Knowledge Management, CIKM '21*, Association for Computing Machinery, New York, NY, USA, 2021, pp. 4680–4684. ISBN 9781450384469. doi:10.1145/3459637.3481981.
- [105] T. Teofili, D. Firmani, N. Koudas, V. Martello, P. Merialdo and D. Srivastava, Effective Explanations for Entity Resolution Models, arXiv, 2022. doi:10.48550/ARXIV.2203.12978. <https://arxiv.org/abs/2203.12978>.
- [106] V. Di Cicco, D. Firmani, N. Koudas, P. Merialdo and D. Srivastava, Interpreting Deep Learning Models for Entity Resolution: An Experience Report Using LIME, in *aiDM '19*, Association for Computing Machinery, New York, NY, USA, 2019. ISBN 9781450368025. doi:10.1145/3329859.3329878.
- [107] X. Mao, W. Wang, Y. Wu and M. Lan, LightEA: A Scalable, Robust, and Interpretable Entity Alignment Framework via Three-view Label Propagation, arXiv, 2022. doi:10.48550/ARXIV.2210.10436. <https://arxiv.org/abs/2210.10436>.
- [108] Z. Yao, C. Li, T. Dong, X. Lv, J. Yu, L. Hou, J. Li, Y. Zhang and Z. Dai, Interpretable and Low-Resource Entity Matching via Decoupling Feature Learning from Decision Making, *CoRR abs/2106.04174* (2021). <https://arxiv.org/abs/2106.04174>.
- [109] T. Deng, L. Hou and Z. Han, Keys as Features for Graph Entity Matching, in: *2020 IEEE 36th International Conference on Data Engineering (ICDE)*, 2020, pp. 1974–1977. doi:10.1109/ICDE48307.2020.00217.
- [110] W. Zhang, S. Deng, H. Wang, Q. Chen, W. Zhang and H. Chen, XTransE: Explainable Knowledge Graph Embedding for Link Prediction with Lifestyles in e-Commerce, in: *Semantic Technology*, X. Wang, F.A. Lisi, G. Xiao and E. Botoeva, eds, Springer Singapore, Singapore, 2020, pp. 78–87.
- [111] J. Stadelmaier and S. Padó, Modeling Paths for Explainable Knowledge Base Completion, in: *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, Association for Computational Linguistics, Florence, Italy, 2019, pp. 147–157. doi:10.18653/v1/W19-4816. <https://aclanthology.org/W19-4816>.
- [112] W. Zhang, B. Paudel, W. Zhang, A. Bernstein and H. Chen, Interaction Embeddings for Prediction and Explanation in Knowledge Graphs, in *WSDM '19*, Association for Computing Machinery, New York, NY, USA, 2019, pp. 96–104. ISBN 9781450359405. doi:10.1145/3289600.3291014.
- [113] U. Zulaika, A. Almeida and D. López-de-Ipiña, Influence Functions for Interpretable link prediction in Knowledge Graphs for Intelligent Environments, in: *2022 7th International Conference on Smart and Sustainable Technologies (SpliTech)*, 2022, pp. 1–7. doi:10.23919/SpliTech55088.2022.9854264.
- [114] W. Jiang, Y. Fu, H. Zhao, J. Wan and S. Pu, Graph Intention Neural Network for Knowledge Graph Reasoning, in: *2022 International Joint Conference on Neural Networks (IJCNN)*, 2022, pp. 1–8. doi:10.1109/IJCNN55064.2022.9892730.
- [115] M.K. Islam, S. Aridhi and M. Smail-Tabbone, Negative sampling and rule mining for explainable link prediction in knowledge graphs, *Knowledge-Based Systems* **250** (2022), 109083. doi:<https://doi.org/10.1016/j.knosys.2022.109083>. <https://www.sciencedirect.com/science/article/pii/S0950705122005342>.
- [116] C. d'Amato, P. Masella and N. Fanizzi, An Approach Based on Semantic Similarity to Explaining Link Predictions on Knowledge Graphs, in: *IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology, WI-IAT '21*, Association for Computing Machinery, New York, NY, USA, 2022, pp. 170–177. ISBN 9781450391153. doi:10.1145/3486622.3493956.
- [117] B.Y. Lin, D. Lee, M. Shen, R. Moreno, X. Huang, P. Shiralkar and X. Ren, TriggerNER: Learning with Entity Triggers as Explanations for Named Entity Recognition, *CoRR abs/2004.07493* (2020). <https://arxiv.org/abs/2004.07493>.
- [118] D. Lee, R.K. Selvam, S.M. Sarwar, B.Y. Lin, M. Agarwal, F. Morstatter, J. Pujara, E. Boschee, J. Allan and X. Ren, AutoTriggER: Named Entity Recognition with Auxiliary Trigger Extraction, *CoRR abs/2109.04726* (2021). <https://arxiv.org/abs/2109.04726>.
- [119] H. Ouchi, J. Suzuki, S. Kobayashi, S. Yokoi, T. Kuribayashi, R. Konno and K. Inui, Instance-Based Learning of Span Representations: A Case Study through Named Entity Recognition, *CoRR abs/2004.14514* (2020). <https://arxiv.org/abs/2004.14514>.
- [120] H. Shahbazi, X. Fern, R. Ghaeini and P. Tadepalli, Relation Extraction with Explanation, in: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, Association for Computational Linguistics, Online, 2020, pp. 6488–6494. doi:10.18653/v1/2020.acl-main.579. <https://aclanthology.org/2020.acl-main.579>.
- [121] A. Albalak, V. Embar, Y. Tuan, L. Getoor and W.Y. Wang, D-REX: Dialogue Relation Extraction with Explanations, *CoRR abs/2109.05126* (2021). <https://arxiv.org/abs/2109.05126>.
- [122] S. Zeng, Y. Wu and B. Chang, SIRE: Separate Intra- and Inter-sentential Reasoning for Document-level Relation Extraction, *CoRR abs/2106.01709* (2021). <https://arxiv.org/abs/2106.01709>.
- [123] Y. Xiao, Z. Zhang, Y. Mao, C. Yang and J. Han, SAIS: Supervising and Augmenting Intermediate Steps for Document-Level Relation Extraction, *CoRR abs/2109.12093* (2021). <https://arxiv.org/abs/2109.12093>.
- [124] W. Zhou, H. Lin, B.Y. Lin, Z. Wang, J. Du, L. Neves and X. Ren, NERO: A Neural Rule Grounding Framework for Label-Efficient Relation Extraction, in: *Proceedings of The Web Conference 2020, WWW '20*, Association for Computing Machinery, New York, NY, USA, 2020, pp. 2166–2176. ISBN 9781450370233. doi:10.1145/3366423.3380282.
- [125] H. Wang, K. Qin, G. Lu, J. Yin, R.Y. Zakari and J.W. Owusu, Document-level relation extraction using evidence reasoning on RST-GRAPH, *Knowledge-Based Systems* **228** (2021), 107274. doi:<https://doi.org/10.1016/j.knosys.2021.107274>. <https://www.sciencedirect.com/science/article/pii/S0950705121005360>.
- [126] H. Kilicoglu, G. Rosemblat, M. Fiszman and D. Shin, Broad-coverage biomedical relation extraction with SemRep, *BMC Bioinformatics* **21**(1) (2020), 188. doi:10.1186/s12859-020-3517-7.
- [127] D. Ru, C. Sun, J. Feng, L. Qiu, H. Zhou, W. Zhang, Y. Yu and L. Li, Learning Logic Rules for Document-level Relation Extraction, *CoRR abs/2111.05407* (2021). <https://arxiv.org/abs/2111.05407>.
- [128] J. Yeo, H. Park, S. Lee, E.W. Lee and S.-w. Hwang, XINA: Explainable Instance Alignment Using Dominance Relationship, *IEEE Transactions on Knowledge and Data Engineering* **32**(2) (2020), 388–401. doi:10.1109/TKDE.2018.2881956.

- [129] R. Singh, V.V. Meduri, A. Elmagarmid, S. Madden, P. Papotti, J.-A. Quiané-Ruiz, A. Solar-Lezama and N. Tang, Synthesizing Entity Matching Rules by Examples, *Proc. VLDB Endow.* **11**(2) (2017), 189–202–. doi:10.14778/3149193.3149199.
- [130] D. Neil, J. Briody, A. Lacoste, A. Sim, P. Creed and A. Saffari, Interpretable Graph Convolutional Neural Networks for Inference on Noisy Knowledge Graphs, *CoRR abs/1812.00279* (2018). <http://arxiv.org/abs/1812.00279>.
- [131] J. Jung, J. Jung and U. Kang, Learning to Walk across Time for Interpretable Temporal Knowledge Graph Completion, in: *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining, KDD '21*, Association for Computing Machinery, New York, NY, USA, 2021, pp. 786–795–. ISBN 9781450383325. doi:10.1145/3447548.3467292.
- [132] Y. Wang, H. Wang, J. He, W. Lu and S. Gao, TAGAT: Type-Aware Graph Attention neTworks for reasoning over knowledge graphs, *Knowledge-Based Systems* **233** (2021), 107500. doi:https://doi.org/10.1016/j.knosys.2021.107500. <https://www.sciencedirect.com/science/article/pii/S0950705121007620>.
- [133] J. Wu, W. Shi, X. Cao, J. Chen, W. Lei, F. Zhang, W. Wu and X. He, DisenKGAT: Knowledge Graph Embedding with Disentangled Graph Attention Network, in: *Proceedings of the 30th ACM International Conference on Information and Knowledge Management, CIKM '21*, Association for Computing Machinery, New York, NY, USA, 2021, pp. 2140–2149–. ISBN 9781450384469. doi:10.1145/3459637.3482424.
- [134] X. Yuan, Q. Lei, S. Yu, C. Xu and Z. Chen, Fine-Grained Relational Learning for Few-Shot Knowledge Graph Completion, *SIGAPP Appl. Comput. Rev.* **22**(3) (2022), 25–38–. doi:10.1145/3570733.3570735.
- [135] J. Wu, S. Mai and H. Hu, Contextual relation embedding and interpretable triplet capsule for inductive relation prediction, *Neurocomputing* **505** (2022), 80–91. doi:https://doi.org/10.1016/j.neucom.2022.07.043. <https://www.sciencedirect.com/science/article/pii/S0925231222008992>.
- [136] R. Bhowmik and G. de Melo, A Joint Framework for Inductive Representation Learning and Explainable Reasoning in Knowledge Graphs, *CoRR abs/2005.00637* (2020). <https://arxiv.org/abs/2005.00637>.
- [137] W. Zhang, S. Deng, M. Chen, L. Wang, Q. Chen, F. Xiong, X. Liu and H. Chen, Knowledge Graph Embedding in E-Commerce Applications: Attentive Reasoning, Explanations, and Transferable Rules, in *IJCKG '21*, Association for Computing Machinery, New York, NY, USA, 2022, pp. 71–79–. ISBN 9781450395656. doi:10.1145/3502223.3502232.
- [138] C. Meilicke, M.W. Chekol, M. Fink and H. Stuckenschmidt, Reinforced Anytime Bottom Up Rule Learning for Knowledge Graph Completion, *CoRR abs/2004.04412* (2020). <https://arxiv.org/abs/2004.04412>.
- [139] Z. Han, P. Chen, Y. Ma and V. Tresp, xERTE: Explainable Reasoning on Temporal Knowledge Graphs for Forecasting Future Links, *CoRR abs/2012.15537* (2020). <https://arxiv.org/abs/2012.15537>.
- [140] A. Sadeghian, M. Armandpour, P. Ding and D.Z. Wang, DRUM: End-To-End Differentiable Rule Mining On Knowledge Graphs, *CoRR abs/1911.00055* (2019). <http://arxiv.org/abs/1911.00055>.
- [141] H. Sun, J. Zhong, Y. Ma, Z. Han and K. He, TimeTraveler: Reinforcement Learning for Temporal Knowledge Graph Forecasting, *CoRR abs/2109.04101* (2021). <https://arxiv.org/abs/2109.04101>.
- [142] S. Ott, C. Meilicke and M. Samwald, SAFRAN: An interpretable, rule-based link prediction method outperforming embedding models, *CoRR abs/2109.08002* (2021). <https://arxiv.org/abs/2109.08002>.
- [143] R. Das, S. Dhuliawala, M. Zaheer, L. Vilnis, I. Durugkar, A. Krishnamurthy, A.J. Smola and A. McCallum, Go for a Walk and Arrive at the Answer: Reasoning Over Paths in Knowledge Bases using Reinforcement Learning, *CoRR abs/1711.05851* (2017). <http://arxiv.org/abs/1711.05851>.
- [144] Y. Liu, Y. Ma, M. Hildebrandt, M. Joblin and V. Tresp, TLogic: Temporal Logical Rules for Explainable Link Forecasting on Temporal Knowledge Graphs, *Proceedings of the AAAI Conference on Artificial Intelligence* **36**(4) (2022), 4120–4127. doi:10.1609/aaai.v36i4.20330. <https://ojs.aaai.org/index.php/AAAI/article/view/20330>.
- [145] W. Xiong, T. Hoang and W.Y. Wang, DeepPath: A Reinforcement Learning Method for Knowledge Graph Reasoning, *CoRR abs/1707.06690* (2017). <http://arxiv.org/abs/1707.06690>.
- [146] P. Wang, K. Agarwal, C. Ham, S. Choudhury and C.K. Reddy, Self-Supervised Learning of Contextual Embeddings for Link Prediction in Heterogeneous Networks, in *WWW '21*, Association for Computing Machinery, New York, NY, USA, 2021, pp. 2946–2957–. ISBN 9781450383127. doi:10.1145/3442381.3450060.
- [147] H. Wang, H. Ren and J. Leskovec, Relational Message Passing for Knowledge Graph Completion, in: *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD '21*, Association for Computing Machinery, New York, NY, USA, 2021, pp. 1697–1707–. ISBN 9781450383325. doi:10.1145/3447548.3467247.
- [148] Z. Du, C. Zhou, J. Yao, T. Tu, L. Cheng, H. Yang, J. Zhou and J. Tang, CogKR: Cognitive Graph for Multi-Hop Knowledge Reasoning, *IEEE Transactions on Knowledge and Data Engineering* **35**(2) (2023), 1283–1295. doi:10.1109/TKDE.2021.3104310.
- [149] C. Lawrence, T. Sztaylor and M. Niepert, Explaining Neural Matrix Factorization with Gradient Rollback, *CoRR abs/2010.05516* (2020). <https://arxiv.org/abs/2010.05516>.
- [150] M. Qu, J. Chen, L.A.C. Xhonneux, Y. Bengio and J. Tang, RNNLogic: Learning Logic Rules for Reasoning on Knowledge Graphs, *CoRR abs/2010.04029* (2020). <https://arxiv.org/abs/2010.04029>.
- [151] C. Fu, T. Chen, M. Qu, W. Jin and X. Ren, Collaborative Policy Learning for Open Knowledge Graph Reasoning, *CoRR abs/1909.00230* (2019). <http://arxiv.org/abs/1909.00230>.
- [152] Y. Gu, Y. Guan and P. Missier, Efficient Rule Learning with Template Saturation for Knowledge Graph Completion, *CoRR abs/2003.06071* (2020). <https://arxiv.org/abs/2003.06071>.
- [153] G. Niu, B. Li, Y. Zhang and S. Pu, CAKE: A Scalable Commonsense-Aware Framework For Multi-View Knowledge Graph Completion, *arXiv*, 2022. doi:10.48550/ARXIV.2202.13785. <https://arxiv.org/abs/2202.13785>.

- [154] M. Hildebrandt, J.A. Quintero Serna, Y. Ma, M. Ringsquandl, M. Joblin and V. Tresp, Reasoning on Knowledge Graphs with Debate Dynamics, *Proceedings of the AAAI Conference on Artificial Intelligence* **34**(04) (2020), 4123–4131. doi:10.1609/aaai.v34i04.6600. <https://ojs.aaai.org/index.php/AAAI/article/view/6600>.
- [155] Y. Zhang and Q. Yao, Knowledge Graph Reasoning with Relational Directed Graph, *CoRR* **abs/2108.06040** (2021). <https://arxiv.org/abs/2108.06040>.
- [156] T. Ma, S. Lv, L. Huang and S. Hu, HiAM: A Hierarchical Attention based Model for knowledge graph multi-hop reasoning, *Neural Networks* **143** (2021), 261–270. doi:<https://doi.org/10.1016/j.neunet.2021.06.008>. <https://www.sciencedirect.com/science/article/pii/S0893608021002409>.
- [157] Y. Xia, M. Lan, J. Luo, X. Chen and G. Zhou, Iterative rule-guided reasoning over sparse knowledge graphs with deep reinforcement learning, *Information Processing & Management* **59**(5) (2022), 103040. doi:<https://doi.org/10.1016/j.ipm.2022.103040>. <https://www.sciencedirect.com/science/article/pii/S0306457322001492>.
- [158] G. Niu, B. Li, Y. Zhang, Y. Sheng, C. Shi, J. Li and S. Pu, Joint semantics and data-driven path representation for knowledge graph reasoning, *Neurocomputing* **483** (2022), 249–261. doi:<https://doi.org/10.1016/j.neucom.2022.02.011>. <https://www.sciencedirect.com/science/article/pii/S0925231222001515>.
- [159] A. Zhu, D. Ouyang, S. Liang and J. Shao, Step by step: A hierarchical framework for multi-hop knowledge graph reasoning with reinforcement learning, *Knowledge-Based Systems* **248** (2022), 108843. doi:<https://doi.org/10.1016/j.knosys.2022.108843>. <https://www.sciencedirect.com/science/article/pii/S0950705122004026>.
- [160] D. Lei, G. Jiang, X. Gu, K. Sun, Y. Mao and X. Ren, Learning Collaborative Agents with Rule Guidance for Knowledge Graph Reasoning, *CoRR* **abs/2005.00571** (2020). <https://arxiv.org/abs/2005.00571>.
- [161] G. Niu, Y. Zhang, B. Li, P. Cui, S. Liu, J. Li and X. Zhang, Rule-Guided Compositional Representation Learning on Knowledge Graphs, *CoRR* **abs/1911.08935** (2019). <http://arxiv.org/abs/1911.08935>.
- [162] C. Zhang, C.-N. Hsu, Y. Katsis, H.-C. Kim and Y. Vázquez-Baeza, Theoretical Rule-based Knowledge Graph Reasoning by Connectivity Dependency Discovery, in: *2022 International Joint Conference on Neural Networks (IJCNN)*, 2022, pp. 1–9. doi:10.1109/IJCNN55064.2022.9891938.
- [163] P. Betz, C. Meilicke and H. Stuckenschmidt, Supervised Knowledge Aggregation for Knowledge Graph Completion, in: *The Semantic Web*, P. Groth, M.-E. Vidal, F. Suchanek, P. Szekley, P. Kapanipathi, C. Pesquita, H. Skaf-Molli and M. Tamper, eds, Springer International Publishing, Cham, 2022, pp. 74–92.
- [164] T. Rocktäschel and S. Riedel, End-to-end Differentiable Proving, in: *Advances in Neural Information Processing Systems*, Vol. 30, I. Guyon, U.V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan and R. Garnett, eds, Curran Associates, Inc., 2017. <https://proceedings.neurips.cc/paper/2017/file/b2ab001909a8a6f04b51920306046ce5-Paper.pdf>.
- [165] Y. Bai, X. Lv, J. Li, L. Hou, Y. Qu, Z. Dai and F. Xiong, SQUIRE: A Sequence-to-sequence Framework for Multi-hop Knowledge Graph Reasoning, *CoRR* **abs/2201.06206** (2022). <https://arxiv.org/abs/2201.06206>.
- [166] G. Niu and B. Li, Logic and Commonsense-Guided Temporal Knowledge Graph Completion, arXiv, 2022. doi:10.48550/ARXIV.2211.16865. <https://arxiv.org/abs/2211.16865>.
- [167] M.A. Hedderich, J. Fischer, D. Klakow and J. Vreeken, Label-Descriptive Patterns and their Application to Characterizing Classification Errors, *CoRR* **abs/2110.09599** (2021). <https://arxiv.org/abs/2110.09599>.
- [168] A. Ebaid, S. Thirumuruganathan, W.G. Aref, A. Elmagarmid and M. Ouzzani, EXPLAINER: Entity Resolution Explanations, in: *2019 IEEE 35th International Conference on Data Engineering (ICDE)*, 2019, pp. 2000–2003. doi:10.1109/ICDE.2019.00224.
- [169] A. Zupon, M. Alexeeva, M. Valenzuela-Escárcega, A. Nagesh and M. Surdeanu, Lightly-supervised Representation Learning with Global Interpretability, in: *Proceedings of the Third Workshop on Structured Prediction for NLP*, Association for Computational Linguistics, Minneapolis, Minnesota, 2019, pp. 18–28. doi:10.18653/v1/W19-1504. <https://aclanthology.org/W19-1504>.
- [170] N. Ding, X. Wang, Y. Fu, G. Xu, R. Wang, P. Xie, Y. Shen, F. Huang, H. Zheng and R. Zhang, Prototypical Representation Learning for Relation Extraction, *CoRR* **abs/2103.11647** (2021). <https://arxiv.org/abs/2103.11647>.
- [171] D.J.T. Cucala, B.C. Grau, E.V. Kostylev and B. Motik, Explainable GNN-Based Models over Knowledge Graphs, in: *International Conference on Learning Representations*, 2022. <https://openreview.net/forum?id=CrCvGNHArz>.
- [172] P. Pezeshkpour, Y. Tian and S. Singh, Investigating Robustness and Interpretability of Link Prediction via Adversarial Modifications, *CoRR* **abs/1905.00563** (2019). <http://arxiv.org/abs/1905.00563>.
- [173] Q. Xie, X. Ma, Z. Dai and E.H. Hovy, An Interpretable Knowledge Transfer Model for Knowledge Base Completion, *CoRR* **abs/1704.05908** (2017). <http://arxiv.org/abs/1704.05908>.
- [174] H. Lu, H. Hu and X. Lin, DensE: An enhanced non-commutative representation for knowledge graph embedding with adaptive semantic hierarchy, *Neurocomputing* **476** (2022), 115–125. doi:<https://doi.org/10.1016/j.neucom.2021.12.079>. <https://www.sciencedirect.com/science/article/pii/S0925231221019342>.
- [175] A. Bastos, K. Singh, A. Nadgeri, S. Shekarpour, I.O. Mulang and J. Hoffart, HopfE: Knowledge Graph Representation Learning Using Inverse Hopf Fibrations, in: *Proceedings of the 30th ACM International Conference on Information and Knowledge Management, CIKM '21*, Association for Computing Machinery, New York, NY, USA, 2021, pp. 89–99. ISBN 9781450384469. doi:10.1145/3459637.3482263.
- [176] Y. Wang, H. Wang, W. Lu and Y. Yan, METransE: Manifold-like mechanism enhanced embedding for reasoning over knowledge graphs, *Expert Systems with Applications* **209** (2022), 118288. doi:<https://doi.org/10.1016/j.eswa.2022.118288>. <https://www.sciencedirect.com/science/article/pii/S0957417422014245>.
- [177] S. Vadrevu, R. Nagi, J. Xiong and W.-m. Hwu, xER: An Explainable Model for Entity Resolution using an Efficient Solution for the Clique Partitioning Problem, in: *Proceedings of the First Workshop on Trustworthy Natural Language Processing*, Association for Computational Linguistics, Online, 2021, pp. 34–44. doi:10.18653/v1/2021.trustnlp-1.5. <https://aclanthology.org/2021.trustnlp-1.5>.

- [178] Q. Lin, R. Mao, J. Liu, F. Xu and E. Cambria, Fusing topology contexts and logical rules in language models for knowledge graph completion, *Information Fusion* **90** (2023), 253–264. doi:<https://doi.org/10.1016/j.inffus.2022.09.020>. <https://www.sciencedirect.com/science/article/pii/S1566253522001592>.
- [179] F. Yang, Z. Yang and W.W. Cohen, Differentiable Learning of Logical Rules for Knowledge Base Completion, *CoRR* **abs/1702.08367** (2017). <http://arxiv.org/abs/1702.08367>.
- [180] A. Meroño-Peñuela, R. Pernisch, C. Guéret and S. Schlobach, Multi-Domain and Explainable Prediction of Changes in Web Vocabularies, in: *Proceedings of the 11th on Knowledge Capture Conference, K-CAP '21*, Association for Computing Machinery, New York, NY, USA, 2021, pp. 193–200–. ISBN 9781450384575. doi:10.1145/3460210.3493583.
- [181] T.-K. Tran, M.H. Gad-Elrab, D. Stepanova, E. Kharlamov and J. Strötgen, Fast Computation of Explanations for Inconsistency in Large-Scale Knowledge Graphs, in: *Proceedings of The Web Conference 2020, WWW '20*, Association for Computing Machinery, New York, NY, USA, 2020, pp. 2613–2619–. ISBN 9781450370233. doi:10.1145/3366423.3380014.
- [182] M. Kejriwal, R. Shao and P. Szekely, Expert-Guided Entity Extraction Using Expressive Rules, in *SIGIR'19*, Association for Computing Machinery, New York, NY, USA, 2019. ISBN 9781450361729. doi:10.1145/3331184.3331392.
- [183] K. Qian, L. Popa and P. Sen, SystemER: A Human-in-the-Loop System for Explainable Entity Resolution **12**(12) (2019), 1794–1797–. doi:10.14778/3352063.3352068.
- [184] M. Paganelli, P. Sottovia, F. Guerra and Y. Velegrakis, TuneR: Fine Tuning of Rule-Based Entity Matchers, in: *Proceedings of the 28th ACM International Conference on Information and Knowledge Management, CIKM '19*, Association for Computing Machinery, New York, NY, USA, 2019, pp. 2945–2948–. ISBN 9781450369763. doi:10.1145/3357384.3357854.
- [185] Y. He, J. Chen, D. Antonyrajah and I. Horrocks, BERTMap: A BERT-based Ontology Alignment System, *CoRR* **abs/2112.02682** (2021). <https://arxiv.org/abs/2112.02682>.
- [186] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A.N. Gomez, L. Kaiser and I. Polosukhin, Attention Is All You Need, *CoRR* **abs/1706.03762** (2017). <http://arxiv.org/abs/1706.03762>.
- [187] M. Kejriwal, A meta-engine for building domain-specific search engines, *Software Impacts* **7** (2021), 100052. doi:<https://doi.org/10.1016/j.simpa.2020.100052>. <https://www.sciencedirect.com/science/article/pii/S2665963820300439>.
- [188] G. Plumb, D. Molitor and A. Talwalkar, Model agnostic supervised local explanations, in: *Proceedings of the 32nd International Conference on Neural Information Processing Systems, NIPS'18*, Curran Associates Inc., Red Hook, NY, USA, 2018, pp. 2520–2529–.
- [189] D. Bahdanau, K. Cho and Y. Bengio, Neural Machine Translation by Jointly Learning to Align and Translate, in: *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Y. Bengio and Y. LeCun, eds, 2015. <http://arxiv.org/abs/1409.0473>.
- [190] R. Singh, V. Meduri, A. Elmagarmid, S. Madden, P. Papotti, J.-A. Quiané-Ruiz, A. Solar-Lezama and N. Tang, Generating Concise Entity Matching Rules, in: *Proceedings of the 2017 ACM International Conference on Management of Data, SIGMOD '17*, Association for Computing Machinery, New York, NY, USA, 2017, pp. 1635–1638–. ISBN 9781450341974. doi:10.1145/3035918.3058739.
- [191] V. Dignum, *Responsible Artificial Intelligence: How to Develop and Use AI in a Responsible Way*, 1st edn, Springer Publishing Company, Incorporated, 2019. ISBN 3030303705.
- [192] S. Hooker, D. Erhan, P.-J. Kindermans and B. Kim, A Benchmark for Interpretability Methods in Deep Neural Networks, in: *Advances in Neural Information Processing Systems*, Vol. 32, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox and R. Garnett, eds, Curran Associates, Inc., 2019. https://proceedings.neurips.cc/paper_files/paper/2019/file/fe4b855600d0f0cae99daa5c5e5a410-Paper.pdf.
- [193] X. Lv, Y. Cao, L. Hou, J. Li, Z. Liu, Y. Zhang and Z. Dai, Is Multi-Hop Reasoning Really Explainable? Towards Benchmarking Reasoning Interpretability, in: *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Online and Punta Cana, Dominican Republic, 2021, pp. 8899–8911. doi:10.18653/v1/2021.emnlp-main.700. <https://aclanthology.org/2021.emnlp-main.700>.
- [194] M.L. Leavitt and A.S. Morcos, Towards falsifiable interpretability research, *CoRR* **abs/2010.12016** (2020). <https://arxiv.org/abs/2010.12016>.
- [195] M.T. Ribeiro, S. Singh and C. Guestrin, Anchors: High-Precision Model-Agnostic Explanations, *Proceedings of the AAAI Conference on Artificial Intelligence* **32**(1) (2018). doi:10.1609/aaai.v32i1.11491. <https://ojs.aaai.org/index.php/AAAI/article/view/11491>.
- [196] B. Letham, C. Rudin, T.H. McCormick and D. Madigan, Interpretable classifiers using rules and Bayesian analysis: Building a better stroke prediction model, *The Annals of Applied Statistics* **9**(3) (2015), 1350–1371. doi:10.1214/15-AOAS848.
- [197] P. Choudhary, A. Kramer and datascience.com team, datascienceinc/Skater: Enable Interpretability via Rule Extraction(BRL), Zenodo, 2018. doi:10.5281/zenodo.1198885.
- [198] A. Crisan, M. Drouhard, J. Vig and N. Rajani, Interactive Model Cards: A Human-Centered Approach to Model Documentation, in: *2022 ACM Conference on Fairness, Accountability, and Transparency, FAccT '22*, Association for Computing Machinery, New York, NY, USA, 2022, pp. 427–439–. ISBN 9781450393522. doi:10.1145/3531146.3533108.
- [199] L. Asprino, E. Daga, A. Gangemi and P. Mulholland, Knowledge Graph Construction with a Façade: A Unified Method to Access Heterogeneous Data Sources on the Web, *ACM Trans. Internet Technol.* (2022). doi:10.1145/3555312.
- [200] J. Wei, X. Wang, D. Schuurmans, M. Bosma, E.H. Chi, Q. Le and D. Zhou, Chain of Thought Prompting Elicits Reasoning in Large Language Models, *CoRR* **abs/2201.11903** (2022). <https://arxiv.org/abs/2201.11903>.
- [201] S. Hellmann, J. Lehmann, S. Auer and M. Brümmer, Integrating NLP Using Linked Data, in: *The Semantic Web – ISWC 2013*, H. Alani, L. Kagal, A. Fokoue, P. Groth, C. Biemann, J.X. Parreira, L. Aroyo, N. Noy, C. Welty and K. Janowicz, eds, Springer Berlin Heidelberg, Berlin, Heidelberg, 2013, pp. 98–113. ISBN 978-3-642-41338-4.

- [202] C.D. Bonaventura, L. Siciliani, P. Basile, A. Meroño-Peñuela and B. McGillivray, Is Explanation All You Need? An Expert Survey on LLM-generated Explanations for Abusive Language Detection, 2024, Submitted to CLiC-it 2024: Tenth Italian Conference on Computational Linguistics.
- [203] E.F. Kendall and D.L. McGuinness, Ontology engineering, *Synthesis Lectures on the Semantic Web: Theory and Technology* **9**(1) (2019), i–102.
- [204] M.C. Suárez-Figueroa, A. Gómez-Pérez and M. Fernández-López, The NeOn methodology for ontology engineering, in: *Ontology engineering in a networked world*, Springer, 2011, pp. 9–34.
- [205] A. Halfaker and R.S. Geiger, Ores: Lowering barriers with participatory machine learning in wikipedia, *Proceedings of the ACM on Human-Computer Interaction* **4**(CSCW2) (2020), 1–37.
- [206] P. Groth, E. Simperl, M. van Erp and D. Vrandečić, Knowledge Graphs and their Role in the Knowledge Engineering of the 21st Century (Dagstuhl Seminar 22372), *Dagstuhl Reports* **12**(9) (2023), 60–120. doi:10.4230/DagRep.12.9.60.
- [207] M. Kejrival and P. Szekely, Knowledge Graphs for Social Good: An Entity-Centric Search Engine for the Human Trafficking Domain, *IEEE Transactions on Big Data* **8**(3) (2022), 592–606. doi:10.1109/TBDATA.2017.2763164.
- [208] A. Saxena, A. Tripathi and P. Talukdar, Improving Multi-hop Question Answering over Knowledge Graphs using Knowledge Base Embeddings, in: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, Association for Computational Linguistics, Online, 2020, pp. 4498–4507. doi:10.18653/v1/2020.acl-main.412. <https://aclanthology.org/2020.acl-main.412>.
- [209] Y. Lan, G. He, J. Jiang, J. Jiang, W.X. Zhao and J. Wen, Complex Knowledge Base Question Answering: A Survey, *CoRR abs/2108.06688* (2021). <https://arxiv.org/abs/2108.06688>.
- [210] Q. Guo, F. Zhuang, C. Qin, H. Zhu, X. Xie, H. Xiong and Q. He, A Survey on Knowledge Graph-Based Recommender Systems, *IEEE Transactions on Knowledge and Data Engineering* **34**(08) (2022), 3549–3568. doi:10.1109/TKDE.2020.3028705.
- [211] J. Devlin, M. Chang, K. Lee and K. Toutanova, BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, *CoRR abs/1810.04805* (2018). <http://arxiv.org/abs/1810.04805>.
- [212] M. Mitchell, S. Wu, A. Zaldívar, P. Barnes, L. Vasserman, B. Hutchinson, E. Spitzer, I.D. Raji and T. Gebru, Model Cards for Model Reporting, *CoRR abs/1810.03993* (2018). <http://arxiv.org/abs/1810.03993>.
- [213] R. Kumar, P. Raghavan, S. Rajagopalan and A. Tomkins, Extracting large-scale knowledge bases from the web, in: *VLDB*, Vol. 99, Citeseer, 1999, pp. 639–650.
- [214] V.S. Silva, A. Freitas and S. Handschuh, Exploring Knowledge Graphs in an Interpretable Composite Approach for Text Entailment, *Proceedings of the AAAI Conference on Artificial Intelligence* **33**(01) (2019), 7023–7030. doi:10.1609/aaai.v33i01.33017023. <https://ojs.aaai.org/index.php/AAAI/article/view/4682>.
- [215] X. Wang, Y. Ye and A. Gupta, Zero-Shot Recognition via Semantic Embeddings and Knowledge Graphs, in: *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6857–6866. doi:10.1109/CVPR.2018.00717.
- [216] Y. Zhang, H. DING, Z. Shui, Y. Ma, J. Zou, A. Deoras and H. Wang, Language Models as Recommender Systems: Evaluations and Limitations, in: *I (Still) Can't Believe It's Not Better! NeurIPS 2021 Workshop*, 2021. <https://openreview.net/forum?id=hFx3fY7-m9b>.
- [217] C. Peng, F. Xia, M. Naseriparsa and F. Osborne, Knowledge Graphs: Opportunities and Challenges, *Artificial Intelligence Review* (2023), 1–32.
- [218] L. Ehrlinger and W. Wöß, Towards a definition of knowledge graphs., *SEMANTICS (Posters, Demos, SuCESS)* **48**(1–4) (2016), 2.
- [219] I. Balazevic, C. Allen and T.M. Hospedales, Multi-relational Poincaré Graph Embeddings, *CoRR abs/1905.09791* (2019). <http://arxiv.org/abs/1905.09791>.
- [220] X. Wang, H. Ji, C. Shi, B. Wang, Y. Ye, P. Cui and P.S. Yu, Heterogeneous Graph Attention Network, in *WWW '19*, Association for Computing Machinery, New York, NY, USA, 2019, pp. 2022–2032. ISBN 9781450366748. doi:10.1145/3308558.3313562.
- [221] R. Angles, H. Thakkar and D. Tomaszuk, RDF and Property Graphs Interoperability: Status and Issues., *AMW* **2369** (2019).
- [222] T. Berners-Lee, J. Hendler and O. Lassila, The semantic web, *Scientific american* **284**(5) (2001), 34–43.
- [223] N. Noy, Y. Gao, A. Jain, A. Narayanan, A. Patterson and J. Taylor, Industry-Scale Knowledge Graphs: Lessons and Challenges: Five Diverse Technology Companies Show How It's Done, *Queue* **17**(2) (2019), 48–75. doi:10.1145/3329781.3332266.
- [224] A. Zaveri, D. Kontokostas, S. Hellmann, J. Umbrich, M. Färber, F. Bartscherer, C. Menne, A. Rettinger, A. Zaveri, D. Kontokostas, S. Hellmann and J. Umbrich, Linked Data Quality of DBpedia, Freebase, OpenCyc, Wikidata, and YAGO, *Semant. Web* **9**(1) (2018), 77–129. doi:10.3233/SW-170275.
- [225] M.J. Dürst and M. Suignard, Internationalized Resource Identifiers (IRIs), *Request for Comments*, RFC Editor, 2005. doi:10.17487/RFC3987. <https://www.rfc-editor.org/info/rfc3987>.
- [226] R. Gatta, M. Vallati, J. Lenkiewicz, E. Rojas, A. Damiani, L. Sacchi, B. De Bari, A. Dagliati, C. Fernandez-Llatas, M. Montesi et al., Generating and comparing knowledge graphs of medical processes using pMineR, in: *Proceedings of the Knowledge Capture Conference*, 2017, pp. 1–4.
- [227] D. Chaves-Fraga, O. Corcho, F. Yedro, R. Moreno, J. Olías and A. De La Azuela, Systematic Construction of Knowledge Graphs for Research-Performing Organizations, *Information* **13**(12) (2022). doi:10.3390/info13120562. <https://www.mdpi.com/2078-2489/13/12/562>.
- [228] M. Hofer, D. Obraczka, A. Saeedi, H. Köpcke and E. Rahm, Construction of Knowledge Graphs: State and Challenges, 2023.
- [229] D. Nadeau and S. Sekine, A survey of named entity recognition and classification, *Linguisticae Investigationes* **30**(1) (2007), 3–26.
- [230] J. Li, A. Sun, J. Han and C. Li, A Survey on Deep Learning for Named Entity Recognition, *CoRR abs/1812.09449* (2018). <http://arxiv.org/abs/1812.09449>.

- [231] J.R. Finkel, T. Grenager and C. Manning, Incorporating Non-local Information into Information Extraction Systems by Gibbs Sampling, in: *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL'05)*, Association for Computational Linguistics, Ann Arbor, Michigan, 2005, pp. 363–370. doi:10.3115/1219840.1219885. <https://aclanthology.org/P05-1045>.
- [232] E.F.T.K. Sang and F.D. Meulder, Introduction to the CoNLL-2003 Shared Task: Language-Independent Named Entity Recognition, *CoRR* **cs.CL/0306050** (2003). <http://arxiv.org/abs/cs/0306050>.
- [233] J. Li, Y. Sun, R.J. Johnson, D. Sciaky, C.-H. Wei, R. Leaman, A.P. Davis, C.J. Mattingly, T.C. Wieggers and Z. Lu, BioCreative V CDR task corpus: a resource for chemical disease relation extraction, *Database* **2016** (2016), baw068. doi:10.1093/database/baw068.
- [234] A. Smirnova and P. Cudré-Mauroux, Relation Extraction Using Distant Supervision: A Survey, *ACM Comput. Surv.* **51**(5) (2018). doi:10.1145/3241741.
- [235] M. Mintz, S. Bills, R. Snow and D. Jurafsky, Distant supervision for relation extraction without labeled data, in: *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, Association for Computational Linguistics, Suntec, Singapore, 2009, pp. 1003–1011. <https://aclanthology.org/P09-1113>.
- [236] G. Papadakis, D. Skoutas, E. Thanos and T. Palpanas, Blocking and Filtering Techniques for Entity Resolution: A Survey, *ACM Comput. Surv.* **53**(2) (2020). doi:10.1145/3377455.
- [237] W. Shen, J. Wang and J. Han, Entity Linking with a Knowledge Base: Issues, Techniques, and Solutions, *IEEE Transactions on Knowledge and Data Engineering* **27**(2) (2015), 443–460. doi:10.1109/TKDE.2014.2327028.
- [238] M. Zamini, H. Reza and M. Rabiei, A Review of Knowledge Graph Completion, *Information* **13**(8) (2022). doi:10.3390/info13080396. <https://www.mdpi.com/2078-2489/13/8/396>.
- [239] A. Bordes, N. Usunier, A. Garcia-Duran, J. Weston and O. Yakhnenko, Translating Embeddings for Modeling Multi-relational Data, in: *Advances in Neural Information Processing Systems*, Vol. 26, C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani and K.Q. Weinberger, eds, Curran Associates, Inc., 2013. https://proceedings.neurips.cc/paper_files/paper/2013/file/1cecc7a77928ca8133fa24680a88d2f9-Paper.pdf.
- [240] Z. Sun, Z. Deng, J. Nie and J. Tang, RotatE: Knowledge Graph Embedding by Relational Rotation in Complex Space, *CoRR* **abs/1902.10197** (2019). <http://arxiv.org/abs/1902.10197>.
- [241] M. Nickel, V. Tresp and H.-P. Kriegel, A Three-Way Model for Collective Learning on Multi-Relational Data, in: *Proceedings of the 28th International Conference on International Conference on Machine Learning, ICML'11*, Omnipress, Madison, WI, USA, 2011, pp. 809–816-. ISBN 9781450306195.
- [242] T. Trouillon, J. Welbl, S. Riedel, É. Gaussier and G. Bouchard, Complex Embeddings for Simple Link Prediction, *CoRR* **abs/1606.06357** (2016). <http://arxiv.org/abs/1606.06357>.
- [243] R. Socher, D. Chen, C.D. Manning and A. Ng, Reasoning With Neural Tensor Networks for Knowledge Base Completion, in: *Advances in Neural Information Processing Systems*, Vol. 26, C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani and K.Q. Weinberger, eds, Curran Associates, Inc., 2013. https://proceedings.neurips.cc/paper_files/paper/2013/file/b337e84de8752b27eda3a12363109e80-Paper.pdf.
- [244] T. Dettmers, P. Minervini, P. Stenetorp and S. Riedel, Convolutional 2D Knowledge Graph Embeddings, *Proceedings of the AAAI Conference on Artificial Intelligence* **32**(1) (2018). doi:10.1609/aaai.v32i1.11573. <https://ojs.aaai.org/index.php/AAAI/article/view/11573>.
- [245] M. Schlichtkrull, T.N. Kipf, P. Bloem, R. van den Berg, I. Titov and M. Welling, Modeling Relational Data with Graph Convolutional Networks, in: *The Semantic Web*, A. Gangemi, R. Navigli, M.-E. Vidal, P. Hitzler, R. Troncy, L. Hollink, A. Tordai and M. Alam, eds, Springer International Publishing, Cham, 2018, pp. 593–607. ISBN 978-3-319-93417-4.
- [246] D. Nathani, J. Chauhan, C. Sharma and M. Kaul, Learning Attention-based Embeddings for Relation Prediction in Knowledge Graphs, *CoRR* **abs/1906.01195** (2019). <http://arxiv.org/abs/1906.01195>.
- [247] K. Toutanova, D. Chen, P. Pantel, H. Poon, P. Choudhury and M. Gamon, Representing Text for Joint Embedding of Text and Knowledge Bases, in: *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Lisbon, Portugal, 2015, pp. 1499–1509. doi:10.18653/v1/D15-1174. <https://aclanthology.org/D15-1174>.
- [248] A. Bosselut, H. Rashkin, M. Sap, C. Malaviya, A. Celikyilmaz and Y. Choi, COMET: Commonsense Transformers for Automatic Knowledge Graph Construction, *CoRR* **abs/1906.05317** (2019). <http://arxiv.org/abs/1906.05317>.
- [249] S. Hao, B. Tan, K. Tang, B. Ni, H. Zhang, E.P. Xing and Z. Hu, BertNet: Harvesting Knowledge Graphs from Pretrained Language Models, 2022.
- [250] F. Petroni, T. Rocktäschel, P.S.H. Lewis, A. Bakhtin, Y. Wu, A.H. Miller and S. Riedel, Language Models as Knowledge Bases?, *CoRR* **abs/1909.01066** (2019). <http://arxiv.org/abs/1909.01066>.
- [251] S. Razniewski, A. Yates, N. Kassner and G. Weikum, Language Models As or For Knowledge Bases, *CoRR* **abs/2110.04888** (2021). <https://arxiv.org/abs/2110.04888>.
- [252] A. Hur, N. Janjua and M. Ahmed, A Survey on State-of-the-art Techniques for Knowledge Graphs Construction and Challenges ahead, in: *2021 IEEE Fourth International Conference on Artificial Intelligence and Knowledge Engineering (AIKE)*, 2021, pp. 99–103. doi:10.1109/AIKE52691.2021.00021.
- [253] Z.C. Lipton, The Mythos of Model Interpretability, *CoRR* **abs/1606.03490** (2016). <http://arxiv.org/abs/1606.03490>.
- [254] F. Doshi-Velez and B. Kim, Towards A Rigorous Science of Interpretable Machine Learning, 2017.
- [255] A. Abdul, J. Vermeulen, D. Wang, B.Y. Lim and M. Kankanhalli, Trends and Trajectories for Explainable, Accountable and Intelligible Systems: An HCI Research Agenda, in: *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, CHI '18, Association for Computing Machinery, New York, NY, USA, 2018, pp. 1–18-. ISBN 9781450356206. doi:10.1145/3173574.3174156.
- [256] A. Adadi and M. Berrada, Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI), *IEEE Access* **6** (2018), 52138–52160. doi:10.1109/ACCESS.2018.2870052.

- [257] Y. Pruksachatkun, M. Mcateer and S. Majumdar, *Practicing Trustworthy Machine Learning*, O'Reilly Media, 2023. ISBN 9781098120238. <https://books.google.co.uk/books?id=LAqmEAAQBAJ>.
- [258] B. Goodman and S. Flaxman, European Union Regulations on Algorithmic Decision-Making and a “Right to Explanation”, *AI Magazine* **38**(3) (2017), 50–57. doi:10.1609/aimag.v38i3.2741. <https://ojs.aaai.org/aimagazine/index.php/aimagazine/article/view/2741>.
- [259] G.K. Dziugaite, S. Ben-David and D.M. Roy, Enforcing Interpretability and its Statistical Impacts: Trade-offs between Accuracy and Interpretability, *CoRR abs/2010.13764* (2020). <https://arxiv.org/abs/2010.13764>.
- [260] A. Bell, I. Solano-Kamaiko, O. Nov and J. Stoyanovich, It’s Just Not That Simple: An Empirical Study of the Accuracy-Explainability Trade-off in Machine Learning for Public Policy, in: *2022 ACM Conference on Fairness, Accountability, and Transparency, FAccT ’22*, Association for Computing Machinery, New York, NY, USA, 2022, pp. 248–266–. ISBN 9781450393522. doi:10.1145/3531146.3533090.
- [261] A. Jacovi, A. Marasović, T. Miller and Y. Goldberg, Formalizing Trust in Artificial Intelligence: Prerequisites, Causes and Goals of Human Trust in AI, in: *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, FAccT ’21*, Association for Computing Machinery, New York, NY, USA, 2021, pp. 624–635–. ISBN 9781450383097. doi:10.1145/3442188.3445923.
- [262] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller and W. Samek, On Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise Relevance Propagation, *PLOS ONE* **10**(7) (2015), 1–46. doi:10.1371/journal.pone.0130140.
- [263] T. Schnake, O. Eberle, J. Lederer, S. Nakajima, K.T. Schütt, K. Müller and G. Montavon, XAI for Graphs: Explaining Graph Neural Network Predictions by Identifying Relevant Walks, *CoRR abs/2006.03589* (2020). <https://arxiv.org/abs/2006.03589>.
- [264] A. Ali, T. Schnake, O. Eberle, G. Montavon, K.-R. Müller and L. Wolf, XAI for Transformers: Better Explanations through Conservative Propagation, in: *Proceedings of the 39th International Conference on Machine Learning*, K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu and S. Sabato, eds, Proceedings of Machine Learning Research, Vol. 162, PMLR, 2022, pp. 435–451. <https://proceedings.mlr.press/v162/ali22a.html>.
- [265] R. Ying, D. Bourgeois, J. You, M. Zitnik and J. Leskovec, GNN Explainer: A Tool for Post-hoc Explanation of Graph Neural Networks, *CoRR abs/1903.03894* (2019). <http://arxiv.org/abs/1903.03894>.
- [266] Q. Huang, M. Yamada, Y. Tian, D. Singh, D. Yin and Y. Chang, GraphLIME: Local Interpretable Model Explanations for Graph Neural Networks, *CoRR abs/2001.06216* (2020). <https://arxiv.org/abs/2001.06216>.
- [267] S. Jain and B.C. Wallace, Attention is not Explanation, *CoRR abs/1902.10186* (2019). <http://arxiv.org/abs/1902.10186>.
- [268] J. Bastings and K. Filippova, The elephant in the interpretability room: Why use attention as explanation when we have saliency methods?, in: *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, Association for Computational Linguistics, Online, 2020, pp. 149–155. doi:10.18653/v1/2020.blackboxnlp-1.14. <https://aclanthology.org/2020.blackboxnlp-1.14>.
- [269] R. Luo, L. Sun, Y. Xia, T. Qin, S. Zhang, H. Poon and T.-Y. Liu, BioGPT: generative pre-trained transformer for biomedical text generation and mining, *Briefings in Bioinformatics* **23**(6) (2022), bbac409. doi:10.1093/bib/bbac409.
- [270] J. Raad and C. Cruz, A Survey on Ontology Evaluation Methods, in: *Proceedings of the International Joint Conference on Knowledge Discovery, Knowledge Engineering and Knowledge Management, IC3K 2015*, SCITEPRESS - Science and Technology Publications, Lda, Setubal, PRT, 2015, pp. 179–186–. ISBN 9789897581588. doi:10.5220/0005591001790186.
- [271] T. Castro Ferreira, C. Gardent, N. Ilinykh, C. van der Lee, S. Mille, D. Moussallem and A. Shimorina (eds), Proceedings of the 3rd International Workshop on Natural Language Generation from the Semantic Web (WebNLG+), Association for Computational Linguistics, Dublin, Ireland (Virtual), 2020. <https://aclanthology.org/2020.webnlg-1.0>.
- [272] D. Vrandečić, *Ontology Evaluation*, in: *Handbook on Ontologies*, S. Staab and R. Studer, eds, Springer Berlin Heidelberg, Berlin, Heidelberg, 2009, pp. 293–313. ISBN 978-3-540-92673-3. doi:10.1007/978-3-540-92673-3_13.
- [273] A. Gómez-Pérez, *Ontology Evaluation*, in: *Handbook on Ontologies*, S. Staab and R. Studer, eds, Springer Berlin Heidelberg, Berlin, Heidelberg, 2004, pp. 251–273. ISBN 978-3-540-24750-0. doi:10.1007/978-3-540-24750-0_13.
- [274] H. Liu, Y. Perl and J. Geller, Concept placement using BERT trained by transforming and summarizing biomedical ontology structure, *Journal of Biomedical Informatics* **112** (2020), 103607. doi:<https://doi.org/10.1016/j.jbi.2020.103607>. <https://www.sciencedirect.com/science/article/pii/S1532046420302355>.
- [275] H. Ye, N. Zhang, S. Deng, X. Chen, H. Chen, F. Xiong, X. Chen and H. Chen, Ontology-enhanced Prompt-tuning for Few-shot Learning, *CoRR abs/2201.11332* (2022). <https://arxiv.org/abs/2201.11332>.
- [276] Y. He, J. Chen, E. Jiménez-Ruiz, H. Dong and I. Horrocks, Language Model Analysis for Ontology Subsumption Inference, 2023.
- [277] J. Chen, Y. He, Y. Geng, E. Jimenez-Ruiz, H. Dong and I. Horrocks, Contextual Semantic Embeddings for Ontology Subsumption Prediction, 2023.
- [278] V.E. V and P.S. Kumar, Ontology Verbalization using Semantic-Refinement, *CoRR abs/1610.09964* (2016). <http://arxiv.org/abs/1610.09964>.
- [279] W. Chen, Y. Cao, F. Feng, X. He and Y. Zhang, Explainable Sparse Knowledge Graph Completion via High-order Graph Reasoning Network, arXiv, 2022. doi:10.48550/ARXIV.2207.07503. <https://arxiv.org/abs/2207.07503>.
- [280] Q. Wang, Y. Hao and J. Cao, ADRL: An attention-based deep reinforcement learning framework for knowledge graph reasoning, *Knowledge-Based Systems* **197** (2020), 105910. doi:<https://doi.org/10.1016/j.knsys.2020.105910>. <https://www.sciencedirect.com/science/article/pii/S0950705120302525>.
- [281] W.X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, Y. Du, C. Yang, Y. Chen, Z. Chen, J. Jiang, R. Ren, Y. Li, X. Tang, Z. Liu, P. Liu, J.-Y. Nie and J.-R. Wen, A Survey of Large Language Models, *arXiv preprint arXiv:2303.18223* (2023).
- [282] S. Wu, O. Orsoy, S. Lu, V. Dabrovolski, M. Dredze, S. Gehrmann, P. Kambadur, D. Rosenberg and G. Mann, BloombergGPT: A Large Language Model for Finance, *arXiv preprint arXiv:2303.17564* (2023).

- [283] K. Valmeekam, A. Olmo, S. Sreedharan and S. Kambhampati, Large Language Models Still Can't Plan (A Benchmark for LLMs on Planning and Reasoning about Change), *arXiv preprint arXiv:2206.10498* (2022).
- [284] G. Mialon, R. Dessì, M. Lomeli, C. Nalmpantis, R. Pasunuru, R. Raileanu, B. Rozière, T. Schick, J. Dwivedi-Yu, A. Celikyilmaz, E. Grave, Y. LeCun and T. Scialom, Augmented Language Models: a Survey, *arXiv preprint arXiv:2302.07842* (2023).
- [285] Q. Dong, L. Li, D. Dai, C. Zheng, Z. Wu, B. Chang, X. Sun, J. Xu and Z. Sui, A Survey for In-context Learning, *arXiv preprint arXiv:2301.00234* (2022).
- [286] C. Zhen, Y. Shang, X. Liu, Y. Li, Y. Chen and D. Zhang, A Survey on Knowledge-Enhanced Pre-trained Language Models, *arXiv preprint arXiv:2212.13428* (2022).
- [287] D. Yin, L. Dong, H. Cheng, X. Liu, K.-W. Chang, F. Wei and J. Gao, A survey of knowledge-intensive nlp with pre-trained language models, *arXiv preprint arXiv:2202.08772* (2022).
- [288] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy and S. Bowman, GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding, in: *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 353–355. doi:10.18653/v1/W18-5446. <https://aclanthology.org/W18-5446>.
- [289] F. Petroni, A. Piktus, A. Fan, P. Lewis, M. Yazdani, N. De Cao, J. Thorne, Y. Jernite, V. Karpukhin, J. Maillard, V. Plachouras, T. Rocktäschel and S. Riedel, KILT: a Benchmark for Knowledge Intensive Language Tasks, in: *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Association for Computational Linguistics, Online, 2021, pp. 2523–2544. doi:10.18653/v1/2021.naacl-main.200. <https://aclanthology.org/2021.naacl-main.200>.
- [290] P. Liang, R. Bommasani, T. Lee, D. Tsipras, D. Soylu, M. Yasunaga, Y. Zhang, D. Narayanan, Y. Wu, A. Kumar, B. Newman, B. Yuan, B. Yan, C. Zhang, C. Cosgrove, C.D. Manning, C. Ré, D. Acosta-Navas, D.A. Hudson, E. Zelikman, E. Durmus, F. Ladhak, F. Rong, H. Ren, H. Yao, J. Wang, K. Santhanam, L. Orr, L. Zheng, M. Yuksekgonul, M. Suzgun, N. Kim, N. Guha, N. Chatterji, O. Khattab, P. Henderson, Q. Huang, R. Chi, S.M. Xie, S. Santurkar, S. Ganguli, T. Hashimoto, T. Icard, T. Zhang, V. Chaudhary, W. Wang, X. Li, Y. Mai, Y. Zhang and Y. Koreeda, Holistic evaluation of language models, *arXiv preprint arXiv:2211.09110* (2022).
- [291] N. Bian, X. Han, L. Sun, H. Lin, Y. Lu and B. He, ChatGPT is a Knowledgeable but Inexperienced Solver: An Investigation of Commonsense Problem in Large Language Models, *arXiv preprint arXiv:2303.16421* (2023).
- [292] S.J. Russell, *Artificial intelligence: a modern approach*, Pearson Education, Inc., 2010.
- [293] W. Zhong, R. Cui, Y. Guo, Y. Liang, S. Lu, Y. Wang, A. Saied, W. Chen and N. Duan, AGIEval: A Human-Centric Benchmark for Evaluating Foundation Models, *arXiv preprint arXiv:2304.06364* (2023).
- [294] K. Meng, D. Bau, A. Andonian and Y. Belinkov, Locating and Editing Factual Associations in GPT, *Advances in Neural Information Processing Systems* **36** (2022).
- [295] C. Shah and R.W. White, *Task-Based Evaluation*, in: *Task Intelligence for Search and Recommendation*, Springer International Publishing, Cham, 2021, pp. 75–98. ISBN 978-3-031-02326-2. doi:10.1007/978-3-031-02326-2_6.
- [296] R. Porzel and R. Malaka, A task-based approach for ontology evaluation.
- [297] M. Ivanovs, R. Kadikis and K. Ozols, Perturbation-based methods for explaining deep neural networks: A survey, *Pattern Recognition Letters* **150** (2021), 228–234. doi:<https://doi.org/10.1016/j.patrec.2021.06.030>. <https://www.sciencedirect.com/science/article/pii/S0167866521002440>.