

Semantic Web Technologies and Bias in Artificial Intelligence: A Systematic Literature Review

Paula Reyero-Lobo ^{a,*}, Enrico Daga ^a, Harith Alani ^a, and Miriam Fernandez ^a

^a *Knowledge Media Institute, The Open University, United Kingdom*

E-mails: paula.reyero-lobo@open.ac.uk, enrico.daga@open.ac.uk, harith.alani@open.ac.uk, miriam.fernandez@open.ac.uk

Abstract. Bias in artificial intelligence (AI) is a critical and timely issue due to its sociological, economic and legal impact, as decisions made by biased algorithms could lead to unfair treatment of specific individuals or groups. Multiple surveys have emerged to provide a multidisciplinary view of bias or to review bias in specific areas such as social sciences, business research, criminal justice, or data mining. Given the ability of Semantic Web (SW) technologies to support multiple AI systems, we review the extent to which semantics can be a "tool" to address bias in different algorithmic scenarios. We provide an in-depth categorisation and analysis of bias assessment, representation, and mitigation approaches that use SW technologies. We discuss their potential in dealing with issues such as representing disparities of specific demographics or reducing data drifts, sparsity, and missing values. We find research works on AI bias that apply semantics mainly in information retrieval, recommendation and natural language processing applications and argue through multiple use cases that semantics can help deal with technical, sociological, and psychological challenges.

Keywords: Bias in AI, Semantics, Semantic Web technologies, Bias assessment, Bias representation, Bias mitigation, Algorithmic fairness

1. Introduction

There is growing awareness of bias and discrimination in AI applications. Users from inactive groups are more at risk of being mistreated on e-commerce platforms such as Amazon or eBay, which is problematic as these often correspond to limited income groups [1]. One of the main challenges in image searching is its limitation to only the sample set of training data [2, 3], which can lead to irrelevant or inaccurate results but, at worst, incorrect associations that reflect and perpetuate the harm done to historically disadvantaged groups [4]. These are just a few enlightening examples of the use cases covered in this survey article. Understandably, the direction of the AI community is shifting to-

wards the pursuit of not only accurate but also ethical AI [5].

One of the main advantages of AI over human intelligence is its ability to process vast amounts of data. Indeed, data plays a fundamental role by which algorithmic decisions can reproduce or even amplify human biases, as these systems are only as good as the data they work with [6]. One of the main challenges of AI is dealing with data limitations, such as incomplete, unrepresentative and erroneous data [7]. In addition, how humans do or do not have access to these systems and how they interact with them are also key bias factors to take into account in AI design.

The vast amount of information available on the Semantic Web (SW) has enormous potential to address the bias problems mentioned above by leveraging the structured formalisation of machine-understandable knowledge to build more realistic and fairer models. There are examples in different domains, such as ma-

*Corresponding author. E-mail: paula.reyero-lobo@open.ac.uk.

chine learning and data mining, natural language processing or social networking and media representation, where the SW, linked data and the web of data have made a significant contribution [8]. For example, we can leverage semantics to control and restrict personal and sensitive data access, support different AI processing tasks such as reasoning, mining, clustering and learning, or extract arguments from natural language text. These systems are not averse to systemic bias, e.g., they may lack information from specific domains or entities that are less popular than others, or information from specific demographics may have more detail depending on the contributor's interest [9]. While we consider the potential bias that SW technologies may have at the data and schema levels, we mainly focus in this article on the contribution that SW technologies can make as a "tool" to address bias in different algorithmic scenarios to promote algorithmic fairness. This analysis is relevant for the SW and AI communities, as bias is gaining attention in different areas, such as computer science, social sciences, philosophy and law [5].

In this study, we provide a review of the contribution of SW technologies to bias in AI. We aim to explain why bias arises and at what level of the AI system, *to better understand harmful behaviours*, and how bias manifests *to understand better whom it affects and how*. This in-depth conceptualisation is crucial due to the lack of consistency between the motivation and the technological solutions proposed to address bias in AI [10], as we need to understand what system behaviours are considered harmful, in what way, to whom and why.

We follow a systematic approach [11] to review the literature and analyse the existing bias solutions that use semantics. Specifically, we focus on the following contributions:

- i) We provide a survey of 34 papers that use semantic-based techniques to address bias in AI.
- ii) We categorise relevant papers according to the type of semantics used, and the type of bias they target.
- iii) We highlight the most common AI application areas under semantic research for bias.
- iv) We identify further challenges in AI bias research for the SW and AI communities.

The rest of the paper is organised as follows. In Section 2, we describe the methodology of the systematic literature review. In Section 3, we define concepts of semantics and bias used in this article. In Section 4, we

report on an analysis of previous works that use semantics to address bias in AI and discuss the main findings, future opportunities and challenges in Section 5. Finally, we provide a conclusion in Section 6.

2. Survey methodology

To provide a thorough literature review, we followed the guidelines of the systematic mapping study research method [11]. Specifically, we address the following research question (RQ):

To what extent can SW technologies be used to address bias in AI?

Two main components constitute this RQ: *semantics* and *bias*. The first aims to investigate the SW community's contribution in methods, evaluation frameworks, or metrics to address bias problems in AI. The second focuses on bias, aiming to assess the types and sources of bias that semantic knowledge can address and the main challenges in AI that semantics methods can help overcome.

The collection of relevant studies is based on extensive keyword-based querying of the two main elements of two popular scholarly databases: Elsevier Scopus and ISI Web of Knowledge (WoS) (Table 1). We complete our search with Microsoft Academic Search, Semantic Scholar and Google Scholar to do snowballing [12].

Table 1

Keywords used to search for relevant works in scholarly databases. TITLE-ABS-KEY refer to the title, abstract and keywords of the paper, respectively. We use the wildcard * to ensure multiple spelling variations are included in the search results.

Search keywords

TITLE-ABS-KEY('bias*' OR 'debias*')

AND

TITLE-ABS-KEY('knowledge graph*' OR 'knowledge base*' OR 'ontology' OR 'ontologies' OR 'ontological representation' OR 'ontological knowledge' OR 'thesaurus' OR 'thesauri' OR 'conceptual semantic*')

We collect significant works according to specific inclusion criteria (IC):

IC1: Papers written in English.

IC2: Studies published in relevant journals between 2010 and 2021.

IC3: Only papers subjected to peer review, which include published journal papers, as part of conference proceedings or workshop, and book chapters.

A list of venues representative of the papers found includes the International Semantic Web Conference (ISWC), the European Semantic Web Conference (ESWC), the World Wide Web Conference (WWW), the International Conference on Information and Knowledge Management (CIKM), the Conference on Artificial Intelligence (AAAI), the International Joint Conference on Artificial Intelligence (IJCAI), and the Conference on Empirical Methods in Natural Language Processing and International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). Two reviewers filtered the papers in four subsequent steps (Figure 1).

In the *Source-based* filter, we select Computer Science, Mathematics, Engineering, Business, Decision Science and Social Sciences as relevant sources when using Scopus. In WoS, we consider all search results and use them to complete our search by adding all non-duplicate articles to our list of related papers.

The *Metadata-based* filter is a paper screening based on title, abstract, publication venue and publication year to discard papers not relevant to our RQ. We consider project proposals and literature reviews in the discussion, but not as *Use Cases* in our analysis. In the case of papers published in more than one venue, we include their latest version.

The *Content-based* filter consists of a paper screening based on the introduction, conclusion, or full-text, especially in unclear studies. This research paper aims to investigate the SW technologies used in solutions for bias coming from the use or development of AI systems. Therefore, we exclude papers that lack an AI system or the use of SW technologies. For example, some works lack evidence of improving bias (e.g. there is no experiment or vision on how to address bias in recommendation systems [13]). Others do not use semantics in their solution (e.g., they use knowledge graph embeddings but addresses bias using disjoint test classes [14]).

The *Snowballing* process concludes our search by including additional studies from paper citations when reading the filtered papers in more detail. Out of the identified 58 relevant works, 34 are examples of *Use Cases* that use semantics to address bias. The remaining 24 are surveys, position papers, and research works addressing bias within semantic resources, which are relevant to the discussion.

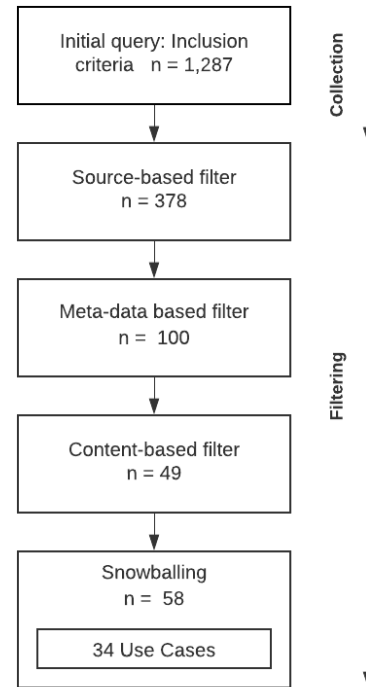


Fig. 1. Filtering relevant works of semantics to address bias in AI.

3. Semantics to address bias in AI

This section presents the definitions for semantics and the conceptualisations of bias in AI used for this paper's analysis.

3.1. Semantics

There are various SW technologies (e.g. taxonomies, thesauri, ontologies, or knowledge graphs). This section defines the specific semantic resources that appear in the surveyed papers to support a better understanding.

Lexical resources are representations of general language in a relational structure [15]. They have standard structured relationships and other properties for each concept, such as related and alternative terms. An example of this type of resource is WordNet [16], a commonly known lexical database of English.

An ontology defines a set of classes, attributes and relationships that model a knowledge domain with varying levels of expressivity [17]. For example, these formal and explicit specifications can be of a shared conceptualisation of meteorological variables (e.g. temperature, precipitation, visibility) to capture the

1 domain of weather forecasting [18], the feelings and
 2 emotions conveyed by visual features to capture the
 3 psychology of human affect [19], or the reasoning
 4 steps of tasks involving problem solving in specific
 5 domains (e.g. writing a risk assessment in industrial,
 6 insurance, health or environmental domains) [20]. In
 7 particular, these forms of *knowledge representation*
 8 can capture terms or statements about the real world
 9 at different levels of domain specialisation [21]. This
 10 scope gives rise to foundational or top-level ontolo-
 11 gies, general or core-reference ontologies. The most
 12 commonly used in the surveyed papers are domain on-
 13 tologies (e.g. for the travel and tourism domain [22]),
 14 and application ontologies (e.g. for the extraction of
 15 information from a weather forecast written in natural
 16 language [18]).

17 A knowledge graph (KG) is a graph of data intended
 18 to accumulate and convey knowledge of the real world,
 19 whose nodes represent entities of interest and whose
 20 edges represent potentially different relations between
 21 these entities [23]. This form of data representation as
 22 a graph represents concepts, classes, properties, rela-
 23 tionships and entity descriptions. ConceptNet [24] is
 24 an example of a KG of 1.6 million assertions of com-
 25 monsense knowledge. A statement such as "cooking
 26 food can be fun" can be represented in the graph as
 27 <cook food> <capableOf> <be fun>. Some common
 28 applications of KGs include recommendation engines,
 29 question answering, or enterprise knowledge manage-
 30 ment [25]. Examples of popular open-source KGs are
 31 DBpedia [26] and Wikidata [27].

32 Finally, Linked Open Data (LOD) refers to a set
 33 of best practises for publishing and connecting struc-
 34 tured data on the Web [28]. LOD relies on documents
 35 containing data in RDF (Resource Description Frame-
 36 work) format to make links between arbitrary *things in*
 37 *the world* (i.e. typed hyperlinks to the related entities
 38 in other data sources). Therefore, instead of navigat-
 39 ing between web pages, Linked Data browsers allow
 40 users to navigate between data sources connected by
 41 specific entities. For example, the LOD platform of the
 42 Open University allows users to navigate the Univer-
 43 sity's content (e.g. courses or scholarly publications)
 44 and establish connections with other educational insti-
 45 tutions [29]. Typically, KGs are published following
 46 the linked data principles.

47 This survey aims to capture the role SW technolo-
 48 gies such as those defined above play in addressing
 49 bias in AI. Specifically, we find research works that use
 50 semantics for three high-level tasks:
 51

1 *Assessing bias* Semantics can uncover bias. As an
 2 example, the representation of user-item interactions
 3 as a graph is used in [1] to discover disparities in the
 4 quality of recommendations in user groups with less
 5 historical data. In algorithmic scenarios where infor-
 6 mation about the protected groups that are vulnerable
 7 to being treated in a discriminatory way, it is of great
 8 value to uncover inequalities through the observable
 9 user properties.

10 *Representing bias* Semantics can capture bias to
 11 make it explicit and raise awareness of its implications.
 12 The use of semantic representations can help to include
 13 information about underrepresented groups in the data
 14 (e.g., due to lack of linguistic coverage in a dataset for
 15 visual sentiment prediction [19]). Documenting con-
 16 sistent errors in black-box models [30] or human when
 17 using such systems [31] can help prevent them and
 18 take action.

19 *Mitigating bias* Semantics can reduce the negative
 20 impact of bias in AI systems. Therefore, we investigate
 21 the combination of SW technologies with methods to
 22 mitigate bias. Bias mitigation methods have generally
 23 been divided into three groups [32, 33]: those focusing
 24 on changing the training data [2, 3, 34–42], the learn-
 25 ing algorithm during the model generation [18, 22, 43–
 26 46], or the model outcomes according to the results
 27 in a holdout dataset which was not involved during
 28 the training phase [1]. Such methods may mitigate un-
 29 desirable associations of specific demographic groups
 30 with hateful connotations. For example, to prevent sen-
 31 tences like "he is gay" or "one of John brothers was ho-
 32 mosexual while the other is a black transgender" from
 33 having high toxicity scores [38].

3.2. Bias in AI

34 Bias can be defined as heterogeneities in data due to
 35 being generated by subgroups of people with their own
 36 characteristics and behaviours [32]. A model learned
 37 from biased data may lead to unfair and inaccurate pre-
 38 dictions. Furthermore, bias can lead to unfairness due
 39 to systematic errors made by algorithms that lead to
 40 adverse or undesired outcomes, for example, of a par-
 41 ticular group that has been historically disadvantaged
 42 [6, 33]. This example of discriminatory behaviour is
 43 particularly concerning since AI systems have proven
 44 to reproduce or even amplify inequalities in society.

45 We aim to help the understanding of bias in AI
 46 through examples that explain why it can occur and
 47 where it comes from, as this knowledge is necessary to
 48
 49
 50
 51

grasp the analysis of the harmful effects and impact of bias in Section 4.

Table 2
Categories of bias attending to the nature of the errors [7].

Type of bias	#Papers	Reference
Cognitive bias	12	[18, 20, 31, 36, 43, 45–51]
Statistical bias	16	[1–3, 22, 34, 35, 39–42, 44, 52–56]
Cultural bias	6	[19, 30, 37, 38, 57, 58]

Table 2 presents the surveyed papers attending to the possible nature of the errors [7]. From a psychological perspective, systematic errors can occur due to the way humans process and interpret information and constitute a *cognitive* bias, which has shown to affect in all decision-making steps [20, 31, 49]. In web search, this can lead to impaired judgement due to the human's heuristic way of processing information [43, 47] and, in more severe cases, to group polarisation [48, 51]. In machine learning (ML), the absence of the context about the domain of the text has shown to leave annotators in an indecisive state, so that their annotations incorrectly shift towards the most frequent sense of a word [50]. The use of subjective text and opinions or any data arising from human interpretation can also have challenging impacts if used to develop AI applications [18, 36, 45, 46].

From a statistical point of view, systematic deviations of the, possibly unknown, real distributions of the variables represented in the data can lead to inaccurate estimations and constitute a *statistical* bias. For example, representation disparities in the data of the users [1], items [41, 53, 54], or their recorded interactions [34] can compromise the quality and fairness of RS. Searching for information based only on the distributions of a specific dataset can lead to irrelevant results or results biased to other meanings of the words used in the query [3, 22]. Similarly, the use of small, domain-specific datasets for training black-box models can lead to undesired behaviours, such as missing the image objects needed to provide meaningful captions [35] and answers to a question about the image [40], or retrieve relevant results to a search query [2], or missing the correct words that enable robots to understand the commands given in a sentence [39]. Consequently, predictions based on these datasets may lead to making decisions based on correlations that are unacceptable in specific cases, as these data only provide estimates from limited settings (e.g. randomised controlled trials [55] or specific datasets commonly used as benchmarks [52]). Such data limitation is especially

concerning in clinical research, as there is a risk of exclusion of women and minorities [56].

From a sociological perspective, data may contain existing biases and beliefs that reflect historical and social inequalities, which the AI system may learn. It constitutes what is known as *cultural* or historical bias. The lack of diversity and overrepresentation of commercial music may lead to music recommendation platforms that are biased towards the specific cultures of popular music [37]. The use of a predominant language, such as English, can lead to generalist systems that are not inclusive of other cultures, for example, when retrieving images based on the feelings they evoke [19] or videos showing what a particular action entails [58]. In some cases, generalised beliefs about particular groups of people are reflected in the data and can lead to learning incorrect associations of these groups with undesirable attributes [30, 38].

Table 3
Categories of bias depending on the location in the AI workflow where bias originates [59].

Bias location	due to	Reference
Bias at source	External bias	[30, 38, 48, 51, 57]
	Functional bias	[1, 34, 41, 43, 47, 53, 56]
Bias at collection	Sampling	[19, 37, 58]
	Querying	[3, 22]
Data pre-processing	Annotation	[18, 36, 39, 45, 46, 50]
	Aggregation	[44]
Data analysis	Inference and prediction	[20, 31, 49]

This categorisation is crucial to reveal how semantics can help with data (i.e. statistical, cultural) and user-dependent (i.e. cognitive) biases. Another aspect to consider is where in the AI workflow these biases originate, so we rely on the framework in [59] for the domain of social media analysis, as it closely reflects general practises in AI.

We use the examples in the surveyed papers shown in Table 3 to discuss this aspect. The first critical point of bias is at the data origin or source since any bias existing at the input of an AI system will appear at least in the same way at the output (i.e. "garbage in - garbage out" principle). This is the bias origin most predominant in the surveyed papers, particularly due to *external* or *functional* factors. The first concerns factors outside the AI system that can influence the reliability and representativeness of the data. For example, the prejudice against specific demographic groups [30, 38], or context of a specific political affiliation [57] or com-

munity views about particular topics [48, 51] may be reflected in the dataset and limit the generalisability of the conclusions that can be drawn from it. The second involves similar limitations due to the design of the AI system. For example, using only purchase data [1], positive feedback [34], or popular items [41, 53] for recommendation affects the data usability. Specific designs can shape and condition users behaviours (e.g. the ranking of search results influences the quality of the information gathered [43, 47]). The heterogeneity of platforms may impede the identification of phenomena analysed on a large scale and also limits the treatment of biases in different study settings [56].

Data collection is the second step in which bias can appear, and examples of this type found in the surveyed papers include *sampling* and *querying* bias. Sampling bias occurs when the data sample is not representative of the whole population (e.g. is only collected from the most popular sources [37] or language [19, 58], so the data collected is not representative of minority groups). Querying bias may emerge due to the lack of expressiveness in the possible query formulations to be able to search for the necessary information (e.g. in an image [3] or information [22] retrieval systems).

Data pre-processing is susceptible to bias, in particular in this study, of *annotation* and *aggregation*. Noisy labels due to poor or missing guidelines compromise manual annotations (e.g. of meteorological analytical data [18], a product review [36], the meaning [50] or link between words in a text [46], or the description of abnormalities in medical images [45]), and frequently lead to the use of small corpora which cannot generalise to novel examples [2, 35, 39, 40, 52, 54]. In domains or problems where the ground truth may not be well defined (e.g. making a medical diagnosis), the use of annotated corpora has limited capacity to ensure that human experts reach a specific level of understanding so that these systems can be applied effectively, efficiently and satisfactorily [55]. An example of bias when transforming the data to infer new facts is found in [44], where bias arises due to an imbalance of the two classes used to infer new sentiment values.

Data analysis is the last step covered by this study that can cause bias, specifically at *inference and prediction* time. For example, issues of this sort may arise when using data as a source of hypotheses rather than a tool to test them [31], or making consistent and predictable mistakes during intelligence activity tasks that draw conclusions from data [20, 49].

4. Description of approaches

Our analysis aims at reaching a better understanding of how bias impacts different AI systems and provides specific methodological examples that can apply to similar problems in future research. This section provides the review of bias assessment, representation and mitigation methodologies that use semantics, which we present following the order in Table 4.

Table 4
Semantic resources used to address bias in AI.

<i>Semantic high-level tasks</i>	<i>Resource</i>	<i>Reference</i>
Bias assessment	Amazon KG	[1]
	Cellosaurus, DrugBank	[56]
	YAGO	[57]
	SentiWordNet	[48]
	<i>prototype</i>	[51]
	FrameNet	[39]
	Wikidata	[52]
	<i>prototype</i>	[55]
	WordNet	[50]
	<i>proprietary</i>	[18]
Bias representation	SentiWordNet	[36]
	<i>medical KG</i>	[45]
	Wikidata	[30]
	<i>prototype</i>	[37]
	<i>MVSO</i>	[19]
	<i>IMAGACT</i>	[58]
	DBpedia	[54]
<i>TIACRITIS</i>	[31]	
Bias mitigation	<i>CBOntology</i>	[20]
	<i>CODM</i>	[49]
	WordNet	[38]
	Wikidata	[43]
	<i>prototype</i>	[47]
	Wikidata	[46]
	DBpedia	[41, 53]
	Freebase	[34]
	ConceptNet	[44]
	<i>prototype</i>	[42]
ConceptNet	[2, 35, 40]	
DBpedia, WebChild	[40]	
ConceptNet, WordNet	[3]	
<i>prototype</i>	[22]	

4.1. Semantics to assess bias

Assessing biases is a fundamental task in analysing and interpreting model behaviours, as it can reveal

1 intrinsic biases that are difficult to detect due to the
2 opaque nature of many AI systems [39]. The follow-
3 ing examples of works are presented as semantics use
4 cases to help with this problem.

5 4.1.1. Bias affecting specific groups of people

6 KGs can help assess recommendation disparities in
7 user groups that are less active [1] (e.g. economically
8 disadvantaged users). This problem constitutes a *pop-*
9 *ulation bias* because it affects groups that are under-
10 represented in the data with respect to the most active
11 users. Therefore, their historic user-interaction data is
12 less visible in the recommendation system. The harm-
13 ful effects of this bias impact the recommendation
14 quality and diversity of the results, posing a fairness
15 problem. A fairness-aware algorithm that leverages en-
16 tities, relationships and paths in KGs is proposed to ex-
17 plicitly model the recommendations in reasoning paths
18 and apply constraints that impose fairness across users.
19 In this case, the Amazon item e-commerce KG of en-
20 tities and relations is used to quantify the richness
21 and evenness of recommendations, revealing dispari-
22 ties of groups with historically less user-item inter-
23 action data. This semantic knowledge is crucial since it
24 is a setting where users do not disclose the personal
25 information required to deal with possible discrimi-
26 natory treatments (i.e. sensitive features such as gen-
27 der, age, or religion). Richness and evenness dispari-
28 ties are measured with each user's number of graph
29 patterns and the relative importance of each pattern
30 across users, respectively. The paper shows how these
31 measures can also be used as fairness constraints to
32 improve the quality and diversity of recommendations
33 for these vulnerable user groups.

34 The problem of under-representation of certain de-
35 mographic groups was assessed in drug exposure stud-
36 ies. Due to the heterogeneity and lack of metadata
37 in the gene expression databases resulting from these
38 medical studies, it is challenging to examine large-
39 scale differences in sex representation of the data. This
40 assessment is essential, as women are 50% more likely
41 to suffer from adverse drug effects. In [56], they were
42 able to assess sex bias in public repositories of bio-
43 logical data by mapping existing and inferred meta-
44 data (using ML models) to existing medical ontolo-
45 gies. Specifically, using Cellosaurus and DrugBank
46 databases to identify cell lines and drugs, respectively.
47 In this way, they could label drug studies from all pub-
48 licly available samples using named entity recognition
49 to identify drug mentions in the metadata and normal-
50 isation to map every instance to all its possible names.
51

1 This analysis served to generate a new resource with
2 unduplicated and normalised data that allows exami-
3 nation across study platforms. As a result, they iden-
4 tified that sex labels are inconsistently reported, with
5 most samples lacking this information. More impor-
6 tantly, they report the existence of sex biases in drug
7 data (e.g., female under-representation in studies of
8 nervous system drugs). This study draws attention to
9 the lack of study and the importance of including sex
10 as a study variable in future analyses.

11 Another use of KGs relevant for this task is to as-
12 sess disparities in the presentation of news reported by
13 different sources (i.e. with different political leaning)
14 [57]. This *media bias* can affect groups or individuals
15 who are part of the story or use these web search sys-
16 tems to build an opinion, since the news are written
17 with the reporter's or media outlet's perception, which
18 can be done, in some cases, partially or unfairly. As
19 a result, bias compromises the reliability of the news
20 source and may raise concerns closely related to the
21 growth of misinformation, polarisation, or online hate.
22 The structure of a KG could help to uncover such dis-
23 parities in reporting between different media outlets,
24 e.g. of specific stakeholders (politicians, political par-
25 ties) advocating or opposing the same issue depend-
26 ing on which source the news appeared. The YAGO
27 KG proved to help extract holders, opinions, and topics
28 and store them, allowing to compare topics and visu-
29 alise biased news. Identifying potentially contradictory
30 information is vital to raise awareness and encourage
31 critical thinking because we, as web users, are exposed
32 to a massive amount of information.

33 4.1.2. Bias affecting individuals

34 The assessment of polarised web search queries is
35 also essential because users are prone to look for infor-
36 mation that reinforces their existing beliefs [48]. Espe-
37 cially when it comes to controversial topics, the *con-*
38 *firmation bias* of web users often conveys views that
39 can lead to a strong division of opinion, affecting in-
40 dividuals and society at large. The impact of bias in
41 this type of user-generated content is closely related to
42 the abovementioned media bias concerns. To prevent
43 these issues, an approach to identifying the sentiment
44 of queries could improve web search systems. This
45 natural language processing (NLP) task (i.e. sentiment
46 analysis) incorporates the support of the SentiWord-
47 Net lexical resource with a two-fold goal. First, it aims
48 to improve the quality of results by including recom-
49 mendations from less popular queries but with similar
50 sentiments. More importantly, providing results from
51

1 queries of opposite sentiment improves the diversity
 2 of opinions and shows the viability of query sentiment
 3 analysis to deal with the problem of bias and polarisa-
 4 tion in web search.

5 Recently, a research project presented a similar ap-
 6 proach to assess confirmation bias and other simi-
 7 lar group phenomena that occur when analysing, vi-
 8 sualising and disseminating information on the Inter-
 9 net (i.e. group polarisation [60] and the belief echo
 10 chamber [61]). The ontology is the primary data struc-
 11 ture for modelling groups and individuals in a net-
 12 work structure [51], and is used as a tool to find simi-
 13 larities and anomalies in the profiles resulting from
 14 the data collected by the system. They present a web
 15 information system to evaluate such effects, integrat-
 16 ing NLP methods into this data structure to identify
 17 themes and sentiment towards them. This architecture
 18 allows tracking individual and group responses to dif-
 19 ferent events (e.g., COVID restrictions, vaccination)
 20 by simply adding new terms to the ontology (e.g., from
 21 specific vaccines such as Pfizer, Moderna or Astra
 22 Zeneca), inferring the sentiment towards them. These
 23 results are promising for integrating AI methods for
 24 natural language understanding (NLU) into distributed
 25 open ecosystems, as this is fundamental to understand-
 26 ing these phenomena in the broader societal domain.

27 4.1.3. Bias affecting AI systems

28 Bias can cause a number of problems affecting AI
 29 systems.

30 *Description of approaches addressing the bias 31 that leads to model overfitting.*

32 Semantics are used to assess inconsistencies in the
 33 predictions of AI systems due to the use of small,
 34 domain-specific corpus for training. *Model overfitting*
 35 can compromise the quality of AI systems and the ex-
 36 tent to which they can fulfil their purpose. It is a sig-
 37 nificant challenge in the AI community due to the po-
 38 tential harms that may arise from using models that are
 39 black-boxes to us, especially when we do not under-
 40 stand why they make specific predictions. In this sur-
 41 vey, we found two use cases that focus on this problem.
 42 In the first example, sentences that entail the same ac-
 43 tion but are different may drift the model towards un-
 44 desired behaviours (i.e., paying attention to the wrong
 45 words in the sentence) [39]. The use of external knowl-
 46 edge from the FrameNet lexical database helped to re-
 47 veal these biased predictions from the mismatch be-
 48 tween the words in a sentence with the highest value
 49 in the attention layer of the model, and the action they
 50
 51

1 should trigger, as captured in this semantic resource.
 2 This analysis revealed patterns that have no theoretical
 3 basis but which the model systematically followed, i.e.
 4 recurrently giving more attention to words that were
 5 not relevant to trigger the action implied by the sen-
 6 tence. This paper also showed how this knowledge
 7 could be included as additional examples in the train-
 8 ing data to make the model more consistent with the
 9 linguistic theory and help it generalise beyond the an-
 10 notated examples of the training corpus.

11 Similar artefacts in datasets are evaluated in pre-
 12 trained masked linguistic models, which are increas-
 13 ingly used in factual knowledge bases to extract in-
 14 formation from a query string [52]. The task con-
 15 sists of using queries such as "Steve Jobs was born
 16 in [MASK]", where *Steve Jobs* is the subject of the
 17 fact and *was born in* a prompt string for the relation
 18 "place of birth", to predict the object placed
 19 as [MASK]. However, their study demonstrates that
 20 many of the successes of *prompt-based* approaches are
 21 due to spurious correlations between similar prompts
 22 (Figure 2). As a result, predictions on completely dif-
 23 ferent datasets are similar, and this is because the
 24 dataset has been over-fitted to specific prompts. The
 25 authors reveal that current *case-based* approaches that
 26 aim to improve performance by providing illustrative
 27 cases mainly succeed in providing a "type guidance".
 28 They reveal that performance is enhanced primarily by
 29 recognising the type of object in the illustrative cases.
 30 Therefore, models can effectively make analogies be-
 31 tween entities of the same type but not predict facts
 32 based on their internal knowledge and the illustrative
 33 cases. This analysis is possible due to the use of the
 34 Wikidata taxonomy to infer the object type of each
 35 relation. It allows us to probe into the behaviours of
 36 these black-box models and better understand the crit-
 37 ical factors underlying their task performance, which
 38 is crucial for building trust in the predictions of these
 39 systems in benchmarks and closed-world studies.

40 Recently, a new framework for automatic decision
 41 making in the medical domain was proposed to meet
 42 the requirements of explainability, robustness, and re-
 43 duced bias in machine learning models [55]. Through
 44 a series of experiments to address three fundamental
 45 challenges in medical research, the authors' reason that
 46 a multi-modal, decentralised and explainable infras-
 47 tructure is needed, where KG can play a crucial role.
 48 In this second use case, a series of arguments are pre-
 49 sented as to why KGs can benefit future human-IA in-
 50 terfaces to be effective in this field (e.g., integrating the
 51 characteristics of different medical data modalities). It

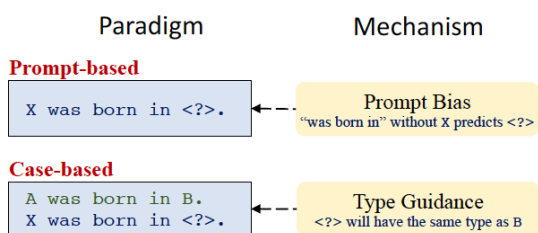


Fig. 2. Example of two dataset artifacts (i.e. *prompt bias*, *type guidance*) that can overfit pre-trained masked language models in factual knowledge extraction tasks [52].

concludes with the basics of *counterfactual graphs*, which store the path from the feature to the changing class to enable the exploration of different counterfactual decision paths (bringing the "human-in-the-loop") and serves as a communication channel with black-box models. This work has a significant impact, as it provides knowledge-based constraints to regularise the training process of deep learning models and the possibility to contest them. This mechanism for opening the black box is critical, as future human-AI interfaces must enable medical experts to understand the causal pathways of automated decision-making systems.

Description of approaches addressing the bias that leads to human annotation errors.

Training corpora also have limitations due to errors in manual annotations that compromise the reliability of the corresponding AI systems. Subjectivity and errors in human annotations constitute a significant problem in developing benchmarks for AI systems, so the assessment of bias in annotation tasks is vital. In [50], background knowledge from the WordNet lexical database is used to support the annotation task where the context (i.e. the set of neighbouring words that provide domain information) is missing. Notably, a comparison of the precision of two lexicographers in a context-agnostic scenario for a word sense disambiguation annotation task, and using WordNet parameters to provide context, revealed that annotations consistently shift towards the most frequent sense of a word in the absence of context. Even though this analysis used few semantic parameters (i.e. conceptual and semantic distance and belonging to the dominant concept), its binding machine versus human annotation study could help demonstrate the importance of context in human annotation tasks.

We found three other examples of assessment of *interpretation bias* in training data. One case study focused on assessing subjectivity in the interpretation of

the analytical variables used to explain weather conditions in forecast texts, as these may vary due to humour, fatigue, or mood [18]. Using this as training data may compromise the truthfulness of the resulting weather prediction systems. They base their approach on the identification of numerical values and properties of different atmospheric variables in texts in order to be able to compare them with observational data. For this, an ontology supports an information extraction model, as it can represent this domain knowledge. Specifically, they developed a proprietary ontology (*AEMIX*) using the Web Ontology Language (OWL) to extract the linguistic information of the important events detected in the text and could reveal the inconsistencies of these texts with the objective information of the interpreted mathematical models.

Similarly, there can be bias when using user product reviews, blog posts and comments to support search engines, recommender systems, and market research applications, due to the subjectivity and ambiguity of this content [36]. As with the previous use case, this may compromise the quality and functionality of NLP systems, especially the degree to which they can detect mixed opinions. They propose a lexical induction approach because mapping subjectivity scores to opinion words in the text can detect review sentiment independently of individual language use. They use SentiWordNet to obtain these sentiment scores, which are used as additional input features for sentiment analysis. Despite the proposed method can only exploit several senses of each word (e.g. the overall score or the value of the first sense) without incorporating semantic relations, it shows the value of ontology-based approaches to avoid human biases arising from the use of machine-learned annotations.

Finally, we have found a good example of how inconsistencies in human-generated data used as ground truth to train automated decision-making systems can compromise the capability and effectiveness of these systems. The use of data-driven neural networks for automated radiology report generation is becoming a critical task in clinical practice [45]. In particular, image captioning approaches trained on medical images and their corresponding reports can significantly improve diagnostic radiology. However, the large volume of images that are a heavy workload for radiologists and, in some cases, lack of experience hinders the generation of these reports. The problem with using previous reports to train automation models is the variability and redundancy between the sentences used to describe the image, especially in describing the nor-

mal regions. For example, as shown in Figure 3, the Blue text corresponds to the description of all normal image elements, while only the Red text indicates the abnormality. Since normal images are already over-represented in the dataset, these deviations and repetitions aggravate the data imbalance and make the generation of sentences to describe normal regions more predominant. As a result, this bias in the data leads to errors where rare but significant abnormalities are not described (Underlined text) and repeated sentences describing the same normal region (Italic text). The use of prior medical domain knowledge captured in a KG showed an improvement in the reports generated, as seen in their quantitative and qualitative results (e.g., the text under "Our" in Figure 3). In particular, the medical KG covering the most common abnormalities and findings can be used as an attention mechanism when exploring new input images. Its framework significantly improves abnormality detection, especially when the occurrence of normal reports dominates the entire dataset. We believe these enlightening examples can motivate future research, as similar data bias issues affect many AI applications.

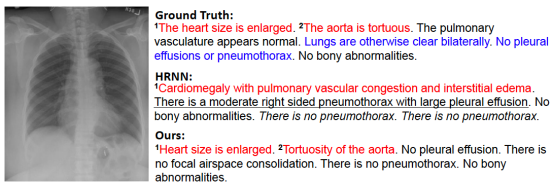


Fig. 3. Example of bias impact in image captioning for automatic radiology report generation when using human-generated data as ground truth for training automated decision-making systems [45].

4.2. Semantics to represent bias

There are some cases in which making the model's systematic preferences and possible biases transparent and expressing them in a human-understandable way may be decisive for providing possible directions of improvement [30]. Therefore, in the following section, we present examples of higher-level semantic tasks for bias representation.

4.2.1. Bias affecting specific groups of people

KGs could help to represent systematic preferences that are consistently applied across the examples used as training data [30]. This is the second example of *population bias*, as these preferences that influence the model, in the same way, are not individual features but

domain categories representing specific groups. Mapping the influential input features to a KG (Wikidata) allowed them to be categorised and described with facts so that groups corresponding to individual entities were captured as *counter-intuitive* rules (e.g. that an Italian origin reduces the value of painters' works). Thus, this additional semantic reasoning capability revealed predictions based on undesirable input data features, such as race or gender, which can be critical for identifying the modification requirements that a model may need to mitigate biases.

There are three other examples of *population bias* where minority groups are disadvantaged by their lack of representation in the data. The first example focuses on the under-representation of minority cultures in music platforms, which leads to a dominance of commercial music and a lack of diversity [37]. Linked web data can be used to represent contextual features of music data (e.g. author biographical information or social connections between artists with similar singing patterns) and create a more relevant navigation space for the cultural background of music. Their prototype is based on a multimodal knowledge base, in which an Open Information Extraction system is used to extract contextual features of music data from the Linked Open Data (LOD). The combination of contextual features together with content features (extracted from audio recordings) can better contextualise the data and reveal non-trivial and deeper relationships between musical entities to lead to more meaningful music discovery and recommendations.

The use of ontologies can be advantageous in improving the representation of minority groups, as shown in the following two examples. In particular, being able to represent differences in the way of expressing and perceiving the affection conveyed by an image across languages can avoid the discrimination of generalist models that only fit a majority language [19]. Language is one of the main characteristics of a culture, so not paying attention to the context of each language can end up damaging entire ethnic groups. In this example, an ontology (Multilingual Visual Sentiment Ontology, *MVSO*) is constructed to represent a training dataset for visual sentiment analysis with a broader scope (including 12 different languages). Using NLP techniques, social media data in these languages and semantic resources (in particular, SentiWordNet and the SentiStrength ontology), they show that including this knowledge in the training datasets improves the degree to which image classification systems can predict sentiment for visual concepts in dif-

ferent languages. This work empirically demonstrates differences in model performance depending on the language used to express the sentiment used in training, as predictions from models trained on a specific language cannot generalise to image data collected in another language. Ensuring diversity in the training data is critical to avoid biased downstream applications to data from the predominant group.

On the other hand, the interpretation of which verbs trigger a given action varies from language to language [58]. Therefore, the development of an ontology with videos to represent different actions can serve as a mechanism to identify groups of verbs that are inclusive of different languages, as the videos serve as a common framework to be annotated in ten different languages. The *IMAGACT* Ontology of Action is a video-based disambiguation framework that can help clustering algorithms not to be specific to a predominant language since they can thus rely on the multilingual lexical features of each action.

4.2.2. Bias affecting AI systems

We present a use case that uses a KG to improve the representation of users' interests beyond the items captured in the training data [54]. This example falls back to the *model overfitting* problem, which in this case can compromise the quality of recommender systems. Their framework incorporates a KG (DBpedia) to expand the user vector representation of a relational graph convolutional network used in the content-based RS. This structure encodes structural and relational information about the neighbouring nodes of the items already part of the training data to provide recommendations consistent with the users' needs. The propagation of relevant knowledge could enhance the performance of the recommender and dialogue systems.

4.2.3. Bias affecting individuals

As final examples, we present three case studies that use semantics to represent consistent and predictable errors that can compromise how data is used and analysed. This *psychological bias* can affect groups developing AI systems to support the search, interpretation, selection and visualisation of information needed to draw conclusions from large masses of data (Intelligence Activity, IA). The first example deals with its impact in the evaluation of evidence, in the search of hypotheses, and argumentation of scientific methods [31]. A domain ontology (*TIACRITIS*) is developed in a collaborative effort to represent all the reasoning steps, probabilistic assessments and assumptions of analysts in data-driven evidence analyses. Similarly,

bias can affect planning, collection, processing and exploitation, analysis and production, dissemination and integration activities [20]. The *CBOntology* is an application ontology that captures the cognitive patterns known to affect these tasks in order to render them explicit and support experts who may experience them. It covers more than 400 classes of such patterns extracted using string, semantic, logical, and topological matching similarities of existing ontologies. These two assistance tools aim to recognise known biases, advise the user to counter them and argue for the need to make biases explicit in AI systems and the experts who use them. The most recent surveyed paper related to this type of bias aims to support and reduce the psychological bias that occurs in decision making under risk and uncertainty [49]. Building on the Core Ontology on Decision Making (*CODM*), the authors extend this knowledge with the decision preferences that are bound in certain circumstances. For this, they use descriptive decision-making theory to extend the ontology with the concepts of intuitive decision-making so that choices made with deliberation or intuition can be explicitly represented to improve understanding of risk preferences and the situation in which they occur. All these works ultimately aim to develop decision support systems that help humans understand their own preferences in order to make better decisions.

4.3. Semantics to mitigate bias

The development of methods to mitigate bias is essential to prevent low-quality results that often impact communities and make them victims of policy injustice, affect their social perceptions, or disadvantage them in other AI application areas [38]. This section presents work examples that leverage semantic knowledge to counteract the possible adverse effects of biased learning.

4.3.1. Bias affecting specific groups of people

Semantic knowledge can be used to mitigate stereotypical perspectives of marginalised groups that are shown to be learned by automated decision-making systems [38]. This is another example of *population bias* because it reflects over-generalised beliefs about specific groups of people that can cause the model to shift towards incorrect predictions. In this example, hate speech detection systems are prone to be overly sensitive to the presence of specific demographic identity terms (e.g. gay, female, black) due to the large amount of hate content that exists against these com-

1 munities. Their technique is based on replacing these
 2 bias-sensitive words with more abstract concepts (e.g.,
 3 gay is a person) to prevent them from being incor-
 4 rectly learned by the model as indicators of hate. To
 5 do so, they use WordNet’s lexical relations to find suit-
 6 able substitution candidates and demonstrate empiri-
 7 cally that systematic deviations towards the hate class
 8 of these terms can be reduced without losing effective-
 9 ness in detecting hate speech. This is an interesting ex-
 10 ample of a bias mitigation technique based on data pre-
 11 processing, as it can reduce this bias in a closed list of
 12 words representing vulnerable groups.

13 4.3.2. Bias affecting individuals

14 KGs have shown advantages in mitigating biases
 15 due to humans’ heuristic way in web searches. We
 16 present two use cases in which the structure of KGs
 17 can help counteract the *confirmation biases* that af-
 18 fect web users. On the one hand, one approach fo-
 19 cuses on investigating the search environment to im-
 20 prove users’ knowledge and attitudes on controversial
 21 topics (in particular, vaccination) [43]. The authors in-
 22 vestigate including factual information extracted from
 23 Wikidata in a knowledge box in the search environ-
 24 ment interface. It showed that users exposed to this in-
 25 formation were significantly more informed, less scep-
 26 tical about vaccination and more critical in discerning
 27 quality information after a simulated web search.

28 Similarly, the advantage of using a KG in the search
 29 interface is addressed in another study that aims to in-
 30 crease the efficiency, quality and user satisfaction with
 31 the information obtained after a web search [47]. To
 32 this end, they developed a KG-based interface pro-
 33 totype using the Open Information Extraction sys-
 34 tem to generate the entity-relationship-entity triplets
 35 of the text. Their qualitative study, based on a post-
 36 experiment evaluation, revealed that the KG interface
 37 helps to reduce the number of times required to view
 38 the source content during exploratory searches with
 39 respect to general hierarchical tree interfaces. These
 40 user-based studies serve to uncover important notions
 41 that shape the use of AI systems.

42 4.3.3. Bias affecting AI systems

43 Bias can cause a number of problems affecting AI
 44 systems.

45 *Description of approaches addressing the bias* 46 *that leads to human annotation errors.*

47 Similar errors can affect the manual annotations of-
 48 ten needed to train AI systems. The *scarcity* of such
 49 data compromises the development of systems to au-
 50 tomate specific tasks. The use case assesses human an-
 51 notators’ errors due to a possible lack of knowledge

1 to provide reliable annotations in extracting informa-
 2 tion from texts [46]. Particularly in cases where differ-
 3 ent mentions corresponding to the same entity have to
 4 be identified (i.e. coreference resolution). A reinforce-
 5 ment learning approach is proposed to address the lack
 6 of examples to train specific neural systems leverag-
 7 ing information from a knowledge base. Specifically,
 8 using Wikidata instances to check the consistency of
 9 facts extracted from the text. Using this information
 10 to tune the model produces better results than other
 11 state-of-the-art methods and paves the way for assist-
 12 ing in the difficult task of obtaining human annotations
 13 needed in many AI systems.

14 *Description of approaches addressing the bias* 15 *that leads to data sparsity.*

16 The rest of the use cases in this section deal with
 17 other limitations regarding how data is used to train
 18 AI systems. One of the problems of recommender ap-
 19 plications is *data sparsity*, as less popular items are
 20 more challenging to deal with and may cause users
 21 to interact only with some of the most popular items
 22 [53]. This can pose a fairness problem because less
 23 popular items are under-represented, and so are the
 24 users that prefer to interact with less popular items.
 25 In this example, a framework is proposed to improve
 26 knowledge-based RS by including the specific seman-
 27 tic properties of the KG. Specifically, the extraction
 28 of each DBpedia property corresponding to user-item
 29 interactions allows computing similarity metrics be-
 30 tween entities that consider each property’s meaning.
 31 These property-specific interactions are included in
 32 the vectors that model the past interactions of each
 33 user to allow making more specific recommendations
 34 (e.g., movies that are related by the actors acting even
 35 if they do not deal with the same topic). This addi-
 36 tional semantic knowledge of user-item interactions
 37 improved recommendations, especially on less popular
 38 data. Specifically, increasing the accuracy of recom-
 39 mended items after discarding the most obvious ones
 40 (serpenticity) and the accuracy of unknown items that
 41 are part of the long tail of the catalogue (novelty).

42 To cope with the large number of features in rec-
 43 ommendation, an entropy-based method is proposed
 44 to obtain only the meaningful historical data of each
 45 user from a KG. The sparse factorisation approach pro-
 46 posed in [41] facilitates the training process by ex-
 47 ploiting a higher level of expressiveness in the fea-
 48 ture embeddings of the items provided by a KG. Facts
 49
 50
 51

and knowledge extracted from DBpedia provide customised recommendation lists, filtering out items with low information gain. Their method allows incorporating the implicit information provided by the KG into the latent space of features, showing an improvement in the quality of results on three benchmark data compared to other state-of-the-art methods. More importantly, their experimental evaluation shows that it improves item diversity, which is critical for measuring popularity bias mitigation.

Description of approaches addressing the bias that leads to missing data.

On the other hand, another limitation imposed by the platforms to collect information for recommendations is the small number of negative samples in the data, as most interactions are positive comments (e.g. clicks, purchases) [34]. This constitutes non-symmetric *missing data*, thus compromising population representation and potentially leading to biased analysis. A KG is proposed to provide informative negative signals to a collaborative filtering algorithm based on matrix factorisation. A negative sampler is constructed using reinforcement learning over Freebase to infer these signals from items related to positive interactions, assuming that these are more likely to be known to the user but were not chosen and, therefore, have a higher probability of being true negatives. Figure 4 is shown as an example. Given that a user (u_1) has watched two movies (i_1, i_2) with the same director (p_1) and genre (p_2), it is more likely that the user knows other movies (e.g. i_4) of the same director but different genre, for which the user has less interest. The reinforcement learning agent over a KG improved the top- K recommendation and preference ranking metrics of seven benchmark methods, which also used KGs, but only to leverage positive signals.

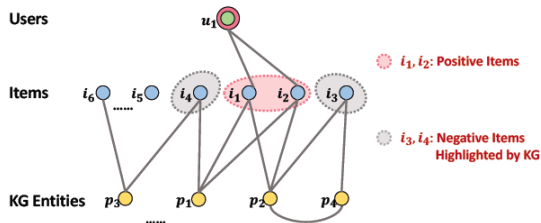


Fig. 4. Example of bias due non-symmetric missing data in a Recommender System (RS) that only collects positive samples [34].

Description of approaches addressing the bias that leads to data imbalance.

Data imbalance can also affect the system's ability to infer new data from existing information. The following example aims to mitigate the bias caused in value propagation methods for sentiment analysis due to the imbalance between positive and negative seeds in the training data [44]. This imbalance causes new inferred values to drift towards the average value gradually. An additional step in the method is proposed to mitigate the bias on the basis that the propagation of values differs depending on the relationship between concepts (e.g. the relationship *isA* has a higher probability of concepts having the same sentiment value than other relations such as *one concept Desires another*). Their method uses a sequential forward search over ConceptNet to select neighbouring concepts with the most relevant type of relationships to propagate new sentiment values. Next, the sentiment value concepts from a manually annotated sentiment dictionary (Affective Norms of English Words, ANEW) are used to align all inferred values with the mean and variance of the concepts that are in the dictionary, assuming that the difference between their inferred and original sentiment values is a shift that occurs due to the imbalance between the initial seeds.

Similarly, bias towards majority samples is a major challenge in current natural language generation (NLG) architectures. As a result, current neural approaches have difficulty generating coherent, grammatically correct text from structured data. The "divide-and-conquer" approach proposed in [42] is based on inducing a hierarchy from a corpus of unlabelled examples using a KG of entity and relation embeddings. The notion of similarity is used to show only the most relevant examples during training to avoid bias due to imbalances in the training data. Specifically, they apply this idea on two datasets containing linked data and textual descriptions (biography paragraphs with semantic mappings to Wikipedia info boxes). Applying similarity of embedded inputs generates effective input-output pairs that consistently outperform competitive baseline approaches. Of particular interest in this study is the partitioning of the dataset according to the semantic and lexical similarity of the entries for training specialised models for each particular similarity group. This general idea can be transferred to other domains to address data sparsity problems (e.g. in image captioning or question answering applications).

Description of approaches addressing the bias that leads to model overfitting.

The following three examples address the problem of *model overfitting* by relying on the use of probabilistic models to generalise better cases that are not included in the training data. In image retrieval [2], relations to concepts in a KG can improve the reasoning power of the model in cases where examples of images with a given caption are not part of the dataset. In particular, the use of the ConceptNet commonsense base can be used to extend the search to images with related captions that are relevant (e.g., the concepts kitchen and restaurant can be informative of chef). This approach can incorporate this rule-based knowledge source and enrich a language model widely used in multimedia-related tasks for NLP. Related concepts, i.e., relations that are relevant in visual space, are included in the object detection function of the model and show improvements in qualitative and quantitative results that are promising for the study of knowledge representation and computer vision.

Second, a similar approach has been applied in an image captioning framework to allow implicit image relationships to be captured in the caption as they may be relevant to describe the image (e.g., if the image shows a "woman standing with her luggage" next to a sign, then it makes sense to speculate that she is waiting for the bus) [35]. The ConceptNet commonsense base helps discover these types of relationships, so a similar strategy is incorporated into the caption generator output to increase the likelihood of latent concepts related to specific objects in the image. This is another example that leverages semantic knowledge to allow the system to generalise beyond the training examples.

Third, the incorporation of external knowledge can help mitigate the errors of systems that respond to questions related to an input image when dealing with answers that did not appear during their training phase or that are not contained within the image scene [40]. In real-world contexts, most techniques fail to address answers that are not within the image content and therefore require external knowledge. For example, Question 1 (Q1) in Figure 5 requires such external knowledge, as responses such as "dog" (an entity that desires the frisbee) cannot be inferred from the image. In this case study, a KG is used to address these limitations, allowing an understanding of the open world scene beyond what is captured in the image. The framework incorporates prior knowledge to guide the alignment process between the feature embeddings of the image-question pair and the corresponding target response. Using a pre-existing subset of DBpedia, ConceptNet and WebChild, this knowl-

edge component is used in which nodes represent the possible set of answers and concepts to enrich the possible relations between them and the possible relations between them. Their quantitative results and their comparison with other current methods support the exploration of knowledge-based systems to overcome overfitting errors.

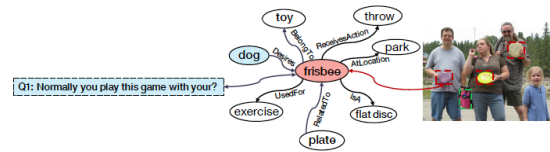


Fig. 5. Example of answer bias in a Visual Question Answering (VQA) system when dealing with concepts not seen during training or in the image scene [40].

Description of approaches addressing the bias that leads to limited expressiveness when using the AI system.

Finally, two case studies have used semantics to mitigate problems in query formulations due to expressiveness limited to a small corpus leading to irrelevant and incorrect results. In image retrieval based on the semantic representation of scene graphs [3], WordNet and ConceptNet can be used to increase the precision of searches that include more complex concepts (e.g. for cases where the system cannot infer that the entities "dog" and "cat" are relevant in the query "animals running on grass"). Their approach introduces a set of rules to find images with fuzzy descriptions and infer the name of the concepts they express, in order to help with more complex searches and enhance semantic and knowledge-based methods for image retrieval processing.

In order to reduce the number of irrelevant results of a web crawler to retrieve social media information, the development of a domain ontology is proposed [22]. Specifically, the Travel&Tour ontology is developed using the Protege-OWL editor to model the specific domain of travel and tourism. The aim is to enrich the content of social media data with the properties and relations of the ontology to improve context-specific searches to take advantage of the domain knowledge provided by this type of data. For example, a search for an expedition-type tourist destination in South America (i.e. *expedition+South+America*) may not return results because the extracted data does not explicitly mention those terms. However, there may be

examples with mentions of related terms relevant to the search (e.g., Amazon river tours offered in the city of Brazil, even though there is no mention of the Amazon being in South America). Although these last two use cases presented only validate their results on a limited set of queries, they are initial works that favour data enrichment to alleviate the lack of expressiveness of the query methods.

5. Discussion

This section highlights the main findings whereby semantics can address bias in AI systems. In addition, it outlines the opportunities for contributions and challenges for further developments in ethical AI.

5.1. Major findings

We summarise some conclusions from the use case analysis according to Table 5 to better elucidate the link of the SW community to bias in AI. We, therefore, highlight the main AI problems and applications where semantics has helped and the tasks for which each semantic technology is most valuable. In particular, we focus our discussion along the following lines:

- i) The use of semantics to address bias in AI is on the rise, and in particular in approaches for mitigating, representing and assessing bias.
- ii) The most researched application areas for biased AI systems using semantics are recommender and search systems and NLP applications.
- iii) KGs are primarily used for bias mitigation, whereas ontologies are mainly used to represent bias, and lexical resources and KGs to assess bias.

Incorporating formal knowledge representations into systems contributes to better generalisation, a fairer balance between bias and accuracy, and more robust methods. There is significant use of KGs to address these *technical challenges*, in particular, to mitigate bias in RS, but also in NLP tasks such as sentiment analysis, image retrieval, image captioning, natural language understanding, natural language generation and visual question answering. We have seen major technical problems are sparsity, missing data, data imbalance, and overfitting due to small and domain-specific training datasets.

The approaches to addressing sociological and psychological challenges have used a more varied range of semantic resources. On the one hand, minority ethnic

groups are often less represented than the general population, which constitutes one of the main *sociological challenges* in AI systems. In this case, approaches generally focused on improving the representation of such groups (often reflected by their language) in different data modalities. In particular, we see how ontologies helped enhance diversity and inclusiveness in multimodal data (i.e. image and video), linked data in music data, and KGs and lexical resources in textual data.

On the other hand, the fact that many AI systems rely on human annotations to learn how to make their future predictions compromises, in many cases, their reliability and truthfulness. Given that human annotations are often liable to subjectivity, interpretation and lack of sufficient knowledge, it constitutes one of the major *psychological challenges* in AI. Previous works use ontologies and lexical resources to assess data annotation problems in NLP tasks (e.g. word sense disambiguation or information extraction from text), whereas KG appears promising for bias mitigation. Therefore, psychological challenges affect the creation of AI systems and how we interact with and use them. We have seen examples using lexical resources and KGs to assess and mitigate confirmation biases in web searches that affect how users interact with and process this content. Finally, ontologies can represent and make explicit the human psychology bias known to occur in computing activities to analyse and draw conclusions from data.

In summary, we can conclude that semantics has contributed to addressing bias that can affect groups, individuals and AI systems, as they pose sociological, psychological and technical challenges.

5.2. Opportunities

We discuss the opportunities for semantics to address bias based on the major findings of this study to elucidate better the connection and future contribution to commonly used methods in the AI bias literature. Specifically, we draw attention to the fact that:

- i) KGs are likely to become increasingly dominant in AI bias research, given their wider scope and potential value in assessing and mitigating bias in various domains.
- ii) Role of semantics in addressing bias related to the collection and annotation of data is likely to become a significant research direction in the near future.

Table 5

Full taxonomy of semantic tasks and bias in AI. Abbrev.: Recommender System (RS), Information Retrieval (IR), Information Extraction (IE), Natural Language Understanding (NLU), Natural Language Generation (NLG), Visual Question Answering (VQA), Text Classification (Text Clf.), Hate Speech detection (HS), Sentiment Analysis (SA), Word Sense Disambiguation (WSD), Music Search and Recommendation (Mus. S/R), Clustering (Clus.), Image Retrieval (Im. R), Image Captioning (Im. C), Image Sentiment Analysis (Im. SA), Intelligence Activity (IA), Content based Filtering (Cnt. F), Collaborative based Filtering (Col. F), Knowledge-based Recommender (K-based R), Scene Graph (SG), Search Engine (SE), Natural Language Processing (NLP), Machine Learning (ML), Computing (Comp.), Linked Open Data (LOD), Lexical Resource (Lexical R), Knowledge Graph (KG).

Type of bias	Bias location	Semantic high-level tasks	AI application	AI technology	SW Technology	Reference
Statistical	Functional	Assessment/Mitigation	RS	K-based R	KG	[1]
		Assessment	Medical Research	ML	Ontology	[56]
		Mitigation	RS	K-based R	KG	[41, 53]
		Mitigation	RS	Col. F	KG	[34]
	Querying	Mitigation	Im. R	SG	KG	[3]
		Mitigation	IR	SE	Ontology	[22]
		Assessment/Mitigation	NLU	NLP	Lexical R	[39]
	Annotation	Assessment	IE	NLP	KG	[52]
		Assessment	Medical Research	ML	KG	[55]
		Representation	RS	Cnt. F	KG	[54]
		Mitigation	NLG	NLP	KG	[42]
		Mitigation	Im. R	NLP	KG	[2]
		Mitigation	Im. C	NLP	KG	[35]
		Mitigation	VQA	NLP	KG	[40]
	Aggregation	Mitigation	SA	NLP	KG	[44]
Cultural	External	Assessment	IR	SE	KG	[57]
		Representation	Text Clf.	NLP	KG	[30]
		Mitigation	HS	NLP	Lexical R	[38]
	Sampling	Representation	Mus. S/R	ML	LOD	[37]
		Representation	Im. SA	NLP	Ontology	[19]
Cognitive	External	Assessment	IR	SE	Lexical R	[48]
		Assessment	IR	SE	Ontology	[51]
	Functional	Mitigation	IR	SE	KG	[43, 47]
		Assessment	WSD	NLP	Lexical R	[50]
	Annotation	Assessment	IE	NLP	Ontology	[18]
		Assessment	SA	NLP	Lexical R	[36]
		Assessment/Mitigation	Im. C	NLP	KG	[45]
	Analysis	Mitigation	IE	NLP	KG	[46]
		Representation	IA	Comp.	Ontology	[20, 31, 49]

iii) Semantic techniques will have a bigger role to play in enabling and enforcing fairness, explainability, and data pre-processing (including data augmentation, enrichment, and correction techniques).

From all the semantic resources that we have analysed in this work, KGs are particularly representative in the last two years (Top-Figure 6). Thus, we expect their use to increase in the coming years. In particular, ConceptNet [2, 3, 35, 44], DBpedia [53, 54] and Wikidata [30, 43] were mostly used to address bias in AI.

Bias mitigation approaches were the most representative of our study period. As seen in the lower part of

Figure 6, semantics can address bias at various stages of the AI workflow and especially the bias coming from the data annotation. The problem lies partly in human annotators' errors when processing information but crucially in the very limitations imposed by using annotated corpora to train AI systems. This is a general practice in AI but leads to systems with a capacity limited to the knowledge captured in a specific dataset.

To address these problems, semantics shows great potential to contribute to several open lines of research on bias in AI. *Fairness metrics* are one of the most well-established approaches to avoid bias and discrimination arising from the data or algorithms used. It is based on measures that evaluate the system's output concerning sensitive or protected attributes that should

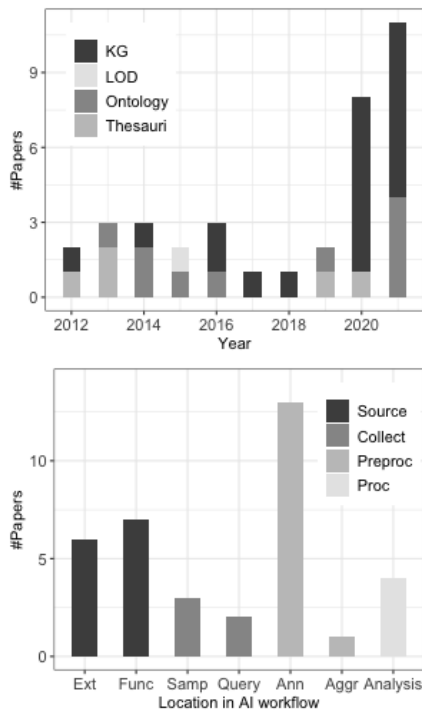


Fig. 6. Semantic web technologies used in the time under scope (Top-Figure). Categories of bias depending on the location in the AI workflow where bias originates (Bottom-Figure). Abbreviations: LOD (Linked Open Data), KG (Knowledge Graph), Source (Bias at source), Collect (Bias at collection), Preproc (Data pre-processing), Proc (Data analysis).

not affect the decision. Semantic knowledge seems promising for retrieving or approximating these characteristics, especially in cases where users have not disclosed this information [1]. Despite the variability and diversity of notation among existing fairness metrics [32], the richness of different perspectives around fairness will hopefully contribute to a better understanding of what fairness is and how to define it in AI systems.

We also emphasise the opportunities of semantics in future *human-AI* interfaces. We encountered several examples where the structure of KGs could enable the addressing of bias in the interaction with AI systems (e.g. when using web search engines [43, 47, 48, 57]). Of particular interest is the vision of knowledge graphs as enablers of interactive and exploration-based explainability techniques [55] as the integration of humans in the training, testing and deployment phases of AI is necessary to bring these systems into real-world contexts.

Finally, we discuss the potential of semantics in the integration of data pre-processing techniques. First, we see a possible intersection of the area of knowledge representation with *data augmentation* techniques to address bias. These techniques rely on increasing samples to deal with unbalanced, unfair distributions or small data sets that may lead to discrimination of specific groups [62]. In this sense, semantics appears useful to re-sample from examples already existing in the data [34] or augment the dataset with new examples to make the model more consistent with the expected behaviours [39, 46]. These approaches seem promising when there are disproportions between different classes or groups.

Furthermore, *data enrichment* techniques address bias by extending the features of instances that are already part of the dataset. In this case, there are several ways to incorporate semantic information into an enriched version of the features. Additional information about features (e.g. properties and relations) can be used as context to improve the generalisation of a specific dataset [37]. However, we must be aware of the noise that may result from the inclusion of uninformative features [63]. Another method is to extract patterns from the graph (e.g. to capture spurious model correlations that are based on sensitive information [30], or properties that enable mining less popular items in a RS [53]). Using KG to represent input features can improve generalisability for models trained with raw input features, as the graph structure gives a further analysis dimension (e.g. for partitioning based on data imbalances using semantic similarity metrics [42]). Finally, probabilistic-based approaches to extend feature vector representations can generalise cases beyond the existing examples in the training dataset, increasing the likelihood of relevant entities given their semantic relationships to a data input. They offer advantages in multiple tasks (image captioning [35, 45], image retrieval [2], visual question answering [40], and recommendation [41, 54]). In summary, contextual enrichment, subgraph pattern mining and probabilistic-based approaches are promising research areas due to the increasing number of cases of individuals, groups, and AI systems still compromised by similar bias problems.

Finally, we discuss the opportunities of *data correction* techniques. Unlike the two previous approaches, these methods modify the data information to account for bias, maintaining the same number of samples and features. Semantic abstraction is a relevant concept in this respect, whereby the use of higher-level concepts

of the information captured in the data can help generalise some dimensions that are not relevant to the task [64]. When data reflects bias and inequalities that the system can learn, such approaches seem worthwhile to retract and reduce the amount of information about specific groups that should not be retained by the model [38].

5.3. Open issues and challenges

In this final part of the discussion, we highlight the challenges and issues we believe future research on AI bias will address.

There is great variability in the evaluation of AI bias-centred research works.

Generally, studies use types of evaluation. User-based evaluation relies on user participation in the system through experimental or observational methods [43, 47, 58]. Besides, a common practice to evaluate the progress of AI systems is using baseline assessments (i.e. comparing approaches using benchmark datasets and specific algorithmic metrics).

The majority of works that address bias evaluate their approaches in downstream implementation. We found works using metrics generally used in recommendation and retrieval applications (ranking scores [2, 34, 40, 54]) and NLP (e.g. textual similarity metrics [35, 42, 45], or general performance in multimodal [19] and text classification tasks [30, 36, 39, 44, 46, 51]). It therefore reveals that there is a need to develop evaluation methods and metrics to assess bias, as an improvement in model performance does not always reflect that the algorithm is not biased. The existence of a possible trade-off between overall performance and bias is an important topic of study in the literature on bias [65] bias. In contrast, only a few previous studies have considered formal definitions of fairness. For example, to ensure the quality and diversity of recommendations in individuals from disadvantaged groups [1] and underrepresented items [41], or to ensure that individuals from specific demographic groups are treated fairly by the system [38].

However, there is an ongoing debate about providing metrics that can be used to benchmark systems addressing bias. In many cases, evaluation frameworks account for demographic information about the individuals or groups affected by AI models. Still, it should take into account various forms of bias in existing models beyond the social categories that are considered as protected attributes by convention [66]. More-

over, these methods for measuring fairness can only reduce discrepancies about the characteristics captured in the data (i.e. the "observed" space) [67]. While these may be relevant for prediction, they may not capture well all the characteristics that served as the basis for decision-making (the "construct" space), leading to the impossibility of a *fair* distribution.

Recommendation. The evaluation of forthcoming semantics-based methods to address bias requires a more critical evaluation that considers potential bias in the context of each particular application.

Secondly, semantic resources cannot be assumed to be free of bias.

Bias can be found in the data used to construct SW technologies.

The concept discussed in [68] of a *polyvocal* and *contextualised* SW draws attention to the fact that these knowledge sources often represent simplified views of the world, in which diverse perspectives may be underrepresented. In this light, the identification, representation and usage of different views or *voices* constitutes one of the main challenges in addressing that SW technologies often reflect the popularity or majority vote. Furthermore, web content is arguably increasingly centralised and asymmetric in terms of the distribution of knowledge and power. Thus, blockchain technologies present themselves as potential next-generation enablers of service exchange and content management [69]. A semantically enriched blockchain software ecosystem based on decentralised applications may be helpful to address bias due to less access of specific demographic groups to these technologies.

There exist several examples in previous work that shed light on the *lack of representation* of specific groups in these technologies. An underrepresentation of less populated countries can occur in manually and semi-automatically created KGs such as Wikidata, as these consequently have a lower number of contributors [70]. Most worryingly, the correlation between coverage and population density is accurate in more developed countries but breaks for the large parts of Asia, Africa, and South America, where their content is drastically underrepresented. Such patterns were consistently found across the different language versions of DBpedia [9].

Demographic bias also propagates in automatic systems to generate KGs. Named entity recognition sys-

tems used for KG construction have shown a systematic exclusion in the detection of entities related to specific demographic categories like gender or ethnicity (e.g. of black female names) [71]. Furthermore, gender disparities can occur in neural relation extraction systems when extracting specific links between entities (e.g. occupation [72]). As a result, bias and prejudice against vulnerable demographic groups propagate into downstream applications. Such harmful associations of specific professions to particular gender, religion, ethnicity and nationality groups (such as men being more likely to be bankers and women to be homemakers) can be found in embeddings extracted from commonly used KGs (i.e. Wikidata, Freebase) [73]. Nevertheless, there are works to address how data representation disparities affect specific demographic groups. One example of a data augmentation method adds new samples to a KG to balance facts that regard specific sensitive attributes (e.g. gender differences in occupations) [74]. This approach effectively mitigates bias in the resulting embeddings from DBpedia and Wikidata. This example stresses the importance of bringing awareness and accounting for the possible bias arising from the application of semantic resources.

Besides the lack of representation, other bias assessed in particular semantic resources is mainly due to low coverage and noise. Disparities in content may exist in different languages. For example, to address *limited coverage* in available general-purpose semantic resources (e.g. to only English), the authors in [75] propose a system to automatically extract lexical FrameNet units using Wikipedia pages in different languages as a reference. However, even when using the same language to model the same domain of knowledge, it is worth bearing in mind that there can be significant disparities between equivalent resources (e.g. as found in the synonym information of four lexical databases [76]).

Statistical methods can serve to estimate the number of facts needed for relations to be representative of the real world [77]. Precisely, a method to calculate a lower bound of different relations could find that at least 46 million facts are missing to draw reasonable conclusions from DBpedia. Many of these missing entities may be due to the lack of type classes covered by DBpedia's ontology and used to automatically extract information from Wikipedia's infoboxes, which leads to only mapping a small subset to the graph [78]. Data descriptive methods appear as potential tools to assess these coverage problems, as shown in the analysis of missing data in specific languages in Wikidata [79].

On the other hand, *noise* is assessed in the annotation, generation, and evaluation of SW technologies. The suitability of the labels given to evaluate semantic-based systems objectively may be compromised, as seen in the differences between expert and crowd-source annotations of natural language summaries generated from a KG [80]. The achievement of link prediction (LP) approaches to automatically extend KGs may be obscured by the existence of inverse and symmetric relations in benchmark datasets. That is, achieving a good performance because certain relations in LP benchmarks tend to occur with others (e.g. the relation `born_in` with `located_in`), or have a default tail answer, as shown empirically in standard benchmarks extracted from Freebase, YAGO, and WordNet [81]. Consequently, the performance of LP methods diminishes tremendously in more realistic settings (e.g. by only removing inverse-duplicate relations from the benchmark dataset [82]). To alleviate the effect of redundant information in the downstream tasks (e.g. information extraction from tabular data to automatically complete KGs [83]), by which only entities well-covered are retained from the table, probabilistic approaches showed useful to return novel facts.

Recommendation. Future semantics-based approaches to address bias in AI should ensure sufficient demographic representation of the people affected by the system and sufficient coverage of the application of use. Additionally, they should use such semantic information in realistic settings that account for noise, mainly due to redundant facts in the captured knowledge.

Consequently, it is imperative to increase transparency and explainability by publishing the source and currency of the data used to generate semantic resources [84], but equally the methods used to construct them. Especially in the enterprise, this information is crucial to ensure the integration of SW technologies in techniques to address bias in AI.

Bias can be found in the methods used to construct SW technologies.

Several factors introducing bias in the development of ontologies have been studied [85]. Specific philosophical views on whether an ontology should represent or interpret reality or its purpose constitute a bias arising from explicit choices. The same is true when capturing insights from competing scientific theories or when economic interests are at stake in deciding which domains deserve more attention. Other factors

1 may propagate bias implicitly, such as specific levels
 2 of granularity, language, or underlying socio-cultural,
 3 political and religious motivations. There are examples
 4 of work that address these limitations (e.g. to define the
 5 scope of an ontology from the literature in a less biased
 6 way towards the selection of particular experts [86], or
 7 to compare the content coverage of the ontology with
 8 a target domain [87]). These examples show the im-
 9 portance of raising awareness of the possible ethical
 10 implications of using these knowledge resources.

11 Similar problems affect the creation of general-
 12 purpose (e.g. DBpedia) and domain-specific KGs (e.g.
 13 GeoNames). Increasingly, KGs are being integrated
 14 into search and recommendation systems to provide
 15 highly personalised content. The problems involved
 16 in creating such personalised KGs can be even more
 17 detrimental than in more conventional methods [88].
 18 Specifically, this representation of users has the risk
 19 of being biased towards specific aspects depending on
 20 the data source used to collect information (e.g. be-
 21 haviours on social networks are different from conver-
 22 sations in forums). In addition, timely events affect the
 23 type of information shared (e.g. in elections or times of
 24 pandemics). This bias can compromise user satisfac-
 25 tion and ultimately aggravate the echo chamber phe-
 26 nomenon. With the evolution and application of se-
 27 mantic technologies in new fields, it is important to be
 28 aware of these new uses' issues.

29 **Recommendation.** The ethical implications of
 30 making particular decisions and selecting partic-
 31 ular sources of information during the develop-
 32 ment of SW technologies require careful consid-
 33 eration when establishing the grounds for their
 34 application in techniques to address bias.

35
 36 To conclude this discussion, we draw on four main
 37 challenges posed by bias and prejudice in AI systems.
 38 First, the need to address the lack of data in unfore-
 39 seen situations (i.e. shifting from controlled to open
 40 environments). We need to develop a world model that
 41 would enable AI for a general-purpose and improve
 42 human-machine communication to use AI as a collab-
 43 orative partner. Finally, we need to establish appropri-
 44 ate trade-offs between conflictive criteria enable these
 45 systems to be applied in a broader range of applica-
 46 tions. Integrating domain knowledge with data-driven
 47 machine learning models in a hybrid approach is key to
 48 address the identified challenges of *environment*, *pur-*
 49 *pose*, *collaboration*, and *governance* [89], so it is pos-
 50 sible to develop ethically sensitive AI methods that
 51 work well in real-world applications. Therefore, while

1 we need a critical analysis before applying semantics
 2 to address bias, the advances in new work and those
 3 seen in this article support research into knowledge-
 4 based reasoning techniques to overcome the pitfalls of
 5 current AI methods.

6. Conclusion

11 This survey article shows the applicability of seman-
 12 tics to address bias in AI. From over a thousand initial
 13 search results, we follow a systematic approach and
 14 present the analysis of 34 use case studies that use for-
 15 mal knowledge representations (i.e. lexical resources,
 16 ontologies, knowledge graphs, or linked data) to as-
 17 sess, represent, or mitigate bias. We provide an ample
 18 understanding and categorisation of bias discussing the
 19 harms associated with bias and the impact it can have
 20 on individuals, groups, and AI systems.

21 Our findings show that semantics has helped in
 22 many AI applications, including information retrieval,
 23 recommender systems, and numerous natural language
 24 processing tasks. Given the increasing use of seman-
 25 tics in recent years, in particular KGs, we conclude
 26 that semantics could especially support fairness, ex-
 27 plainability, and data pre-processing methods to ad-
 28 dress bias in AI. For instance, increasing the sample
 29 size or number of features to provide useful informa-
 30 tion, or manipulating the data samples to account for
 31 bias.

32 We identify further challenges in AI bias research
 33 for the SW and AI communities. These are primarily
 34 the need to develop more robust bias evaluation met-
 35 rics beyond established sensitive information captured
 36 by dataset features. These do not necessarily capture
 37 all the relevant information needed to build fair AI sys-
 38 tems. We also discuss necessary considerations before
 39 applying SW technologies in techniques to address AI
 40 bias, such as underrepresentation of specific demo-
 41 graphic groups, low application coverage, and noise.

42 This paper aims to position the work of the SW com-
 43 munity in the past decade within the context of bias in
 44 AI and provides an analysis of the intersection of both
 45 areas to assist future researchers in identifying and nur-
 46 turing the benefits of these technologies. Bias in AI is
 47 an urgent issue because it compromises the applicabil-
 48 ity of automated systems in society, and the use of se-
 49 mantics has enormous potential to help give meaning
 50 to the data they use.

References

- [1] Z. Fu, Y. Xian, R. Gao, J. Zhao, Q. Huang, Y. Ge, S. Xu, S. Geng, C. Shah, Y. Zhang et al., Fairness-aware explainable recommendation over knowledge graphs, in: *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2020, pp. 69–78.
- [2] R.T. Icarte, J.A. Baier, C. Ruz and A. Soto, How a general-purpose commonsense ontology can improve performance of learning-based image retrieval, *arXiv preprint arXiv:1705.08844* (2017).
- [3] H. Chen, A. Trouve, K.J. Murakami and A. Fukuda, Semantic image retrieval for complex queries using a knowledge parser, *Multimedia Tools and Applications* **77**(9) (2018), 10733–10751.
- [4] T. Simonite, When it comes to gorillas, google photos remains blind, *Wired, January* **13** (2018).
- [5] B. Lepri, N. Oliver and A. Pentland, Ethical machines: The human-centric use of artificial intelligence, *IScience* **24**(3) (2021), 102249.
- [6] S. Barocas and A.D. Selbst, Big data’s disparate impact, *Calif. L. Rev.* **104** (2016), 671.
- [7] R. Baeza-Yates, Bias on the web, *Communications of the ACM* **61**(6) (2018), 54–61.
- [8] F. Gandon, A survey of the first 20 years of research on semantic Web and linked data, *Revue des Sciences et Technologies de l’Information-Série ISI: Ingénierie des Systèmes d’information* (2018).
- [9] K. Janowicz, B. Yan, B. Regalia, R. Zhu and G. Mai, Debiasing Knowledge Graphs: Why Female Presidents are not like Female Popes., in: *International Semantic Web Conference (P&D/Industry/BlueSky)*, 2018.
- [10] S.L. Blodgett, S. Barocas, H. Daumé III and H. Wallach, Language (technology) is power: A critical survey of “bias” in nlp, *arXiv preprint arXiv:2005.14050* (2020).
- [11] P. Brereton, B.A. Kitchenham, D. Budgen, M. Turner and M. Khalil, Lessons from applying the systematic literature review process within the software engineering domain, *Journal of systems and software* **80**(4) (2007), 571–583.
- [12] K. Petersen, S. Vakkalanka and L. Kuzniarz, Guidelines for conducting systematic mapping studies in software engineering: An update, *Information and software technology* **64** (2015), 1–18.
- [13] H.J. Lee and B.-W. Park, How to reduce confirmation bias using linked open data knowledge repository, in: *2020 IEEE International Conference on Big Data and Smart Computing (BigComp)*, IEEE, 2020, pp. 410–416.
- [14] R. Celebi, H. Uyar, E. Yasar, O. Gumus, O. Dikenelli and M. Dumontier, Evaluation of knowledge graph embedding approaches for drug-drug interaction prediction in realistic settings, *BMC bioinformatics* **20**(1) (2019), 1–14.
- [15] F. Richter and M. Sailer, Basic concepts of lexical resource semantics, in: *Esslli*, Citeseer, 2003, pp. 87–143.
- [16] G.A. Miller, *WordNet: An electronic lexical database*, MIT press, 1998.
- [17] T. Gruber, Ontology., *Encyclopedia of database systems* **1** (2009), 1963–1965.
- [18] A.L. Garrido, M.G. Buey, G. Muñoz and J.-L. Casado-Rubio, Information extraction on weather forecasts with semantic technologies, in: *International Conference on Applications of Natural Language to Information Systems*, Springer, 2016, pp. 140–151.
- [19] B. Jou, T. Chen, N. Pappas, M. Redi, M. Topkara and S.-F. Chang, Visual affect around the world: A large-scale multilingual visual sentiment ontology, in: *Proceedings of the 23rd ACM international conference on Multimedia*, 2015, pp. 159–168.
- [20] G. Lortal, P. Capet and A. Bertone, Ontology building for cognitive bias assessment in intelligence, in: *2014 IEEE International Inter-Disciplinary Conference on Cognitive Methods in Situation Awareness and Decision Support (CogSIMA)*, IEEE, 2014, pp. 237–243.
- [21] C. Roussey, F. Pinet, M.A. Kang and O. Corcho, An introduction to ontologies and ontology engineering, in: *Ontologies in Urban development projects*, Springer, 2011, pp. 9–38.
- [22] E. Sedyono, C. Nivak et al., Measuring the performance of ontological based information retrieval from a social media, in: *2014 European Modelling Symposium*, IEEE, 2014, pp. 354–359.
- [23] A. Hogan, E. Blomqvist, M. Cochez, C. d’Amato, G.d. Melo, C. Gutierrez, S. Kirrane, J.E.L. Gayo, R. Navigli, S. Neumaier et al., Knowledge graphs, *Synthesis Lectures on Data, Semantics, and Knowledge* **12**(2) (2021), 1–257.
- [24] H. Liu and P. Singh, ConceptNet—a practical commonsense reasoning tool-kit, *BT technology journal* **22**(4) (2004), 211–226.
- [25] X. Zou, A survey on application of knowledge graph, in: *Journal of Physics: Conference Series*, Vol. 1487, IOP Publishing, 2020, p. 012016.
- [26] S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak and Z. Ives, Dbpedia: A nucleus for a web of open data, in: *The semantic web*, Springer, 2007, pp. 722–735.
- [27] D. Vrandečić and M. Krötzsch, Wikidata: a free collaborative knowledgebase, *Communications of the ACM* **57**(10) (2014), 78–85.
- [28] C. Bizer, T. Heath and T. Berners-Lee, Linked data: The story so far, in: *Semantic services, interoperability and web applications: emerging concepts*, IGI global, 2011, pp. 205–227.
- [29] E. Daga, M. d’Aquin, A. Adamou and S. Brown, The open university linked data—data. open. ac. uk, *Semantic Web* **7**(2) (2016), 183–191.
- [30] A. Nikolov and M. d’Aquin, Uncovering Semantic Bias in Neural Network Models Using a Knowledge Graph, in: *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, 2020, pp. 1175–1184.
- [31] G. Tecuci, D. Schum, D. Marcu and M. Boicu, Recognizing and Countering Biases in Intelligence Analysis with TIA-CRITIS., in: *STIDS*, Citeseer, 2013, pp. 25–32.
- [32] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman and A. Galstyan, A survey on bias and fairness in machine learning, *ACM Computing Surveys (CSUR)* **54**(6) (2021), 1–35.
- [33] E. Ntoutsi, P. Fafalios, U. Gadiraju, V. Iosifidis, W. Nejdl, M.-E. Vidal, S. Ruggieri, F. Turini, S. Papadopoulos, E. Krasnakis et al., Bias in data-driven artificial intelligence systems—An introductory survey, *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* **10**(3) (2020), e1356.
- [34] X. Wang, Y. Xu, X. He, Y. Cao, M. Wang and T.-S. Chua, Reinforced negative sampling over knowledge graph for recommendation, in: *Proceedings of The Web Conference 2020*, 2020, pp. 99–109.

- [35] F. Huang, Z. Li, S. Chen, C. Zhang and H. Ma, Image captioning with internal and external knowledge, in: *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, 2020, pp. 535–544.
- [36] H.-J. Kim and M. Song, An ontology-based approach to sentiment classification of mixed opinions in online restaurant reviews, in: *International Conference on Social Informatics*, Springer, 2013, pp. 95–108.
- [37] G.K. Koduri, Culture-aware approaches to modeling and description of intonation using multimodal data, in: *International Conference on Knowledge Engineering and Knowledge Management*, Springer, 2014, pp. 209–217.
- [38] P. Badjatiya, M. Gupta and V. Varma, Stereotypical bias removal for hate speech detection task using knowledge-based generalizations, in: *The World Wide Web Conference*, 2019, pp. 49–59.
- [39] M. Mensio, E. Bastianelli, I. Tiddi and G. Rizzo, Mitigating bias in deep nets with knowledge bases: The case of natural language understanding for robots, in: *AAAI Spring Symposium: Combining Machine Learning with Knowledge Engineering (1)*, 2020.
- [40] Z. Chen, J. Chen, Y. Geng, J.Z. Pan, Z. Yuan and H. Chen, Zero-Shot Visual Question Answering Using Knowledge Graph, in: *International Semantic Web Conference*, Springer, 2021, pp. 146–162.
- [41] V.W. Anelli, T. Di Noia, E. Di Sciascio, A. Ferrara and A.C.M. Mancino, Sparse feature factorization for recommender systems with knowledge graphs, in: *Fifteenth ACM Conference on Recommender Systems*, 2021, pp. 154–165.
- [42] N. Dethlefs, A. Schoene and H. Cuayáhuitl, A divide-and-conquer approach to neural natural language generation from structured data, *Neurocomputing* **433** (2021), 300–309.
- [43] R. Ludolph, A. Allam, P.J. Schulz et al., Manipulating Google’s knowledge graph box to counter biased information processing during an online search on vaccination: application of a technological debiasing strategy, *Journal of medical Internet research* **18**(6) (2016), e5430.
- [44] C.-E. Wu and R.T.-H. Tsai, Using relation selection to improve value propagation in a conceptnet-based sentiment dictionary, *Knowledge-Based Systems* **69** (2014), 100–107.
- [45] F. Liu, X. Wu, S. Ge, W. Fan and Y. Zou, Exploring and distilling posterior and prior knowledge for radiology report generation, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 13753–13762.
- [46] R. Aralikkatte, H. Lent, A.V. Gonzalez, D. Hershovich, C. Qiu, A. Sandholm, M. Ringgaard and A. Sogaard, Rewarding coreference resolvers for being consistent with world knowledge, *arXiv preprint arXiv:1909.02392* (2019).
- [47] B. Sarrafzadeh, A. Vtyurina, E. Lank and O. Vechtomova, Knowledge graphs versus hierarchies: An analysis of user behaviours and perspectives in information seeking, in: *Proceedings of the 2016 ACM on Conference on Human Information Interaction and Retrieval*, 2016, pp. 91–100.
- [48] S. Chelaru, I.S. Altingovde, S. Siersdorfer and W. Nejdl, Analyzing, detecting, and exploiting sentiment in web queries, *ACM Transactions on the Web (TWEB)* **8**(1) (2013), 1–28.
- [49] E. da Costa Ramos^o, M.L.M. Campos, F. Baião and R. Guizardi^o, Extending the Core Ontology on Decision Making according to Behavioral Economics (2021).
- [50] A. Chatterjee, S. Joshi, P. Bhattacharyya, D. Kanojia and A.K. Meena, A Study of the Sense Annotation Process: Man v/s Machine., in: *GWC 2012 6th International Global Wordnet Conference*, 2012, p. 79.
- [51] M. Pavlíček, T. Filip and P. Sosík, ZREC architecture for textual sentiment analysis (2021).
- [52] B. Cao, H. Lin, X. Han, L. Sun, L. Yan, M. Liao, T. Xue and J. Xu, Knowledgeable or Educated Guess? Revisiting Language Models as Knowledge Bases, *arXiv preprint arXiv:2106.09231* (2021).
- [53] E. Palumbo, D. Monti, G. Rizzo, R. Troncy and E. Baralis, entity2rec: Property-specific knowledge graph embeddings for item recommendation, *Expert Systems with Applications* **151** (2020), 113235.
- [54] Q. Chen, J. Lin, Y. Zhang, M. Ding, Y. Cen, H. Yang and J. Tang, Towards knowledge-based recommender dialog system, *arXiv preprint arXiv:1908.05391* (2019).
- [55] A. Holzinger, B. Malle, A. Saranti and B. Pfeifer, Towards multi-modal causability with Graph Neural Networks enabling information fusion for explainable AI, *Information Fusion* **71** (2021), 28–37.
- [56] E. Flynn, A. Chang and R.B. Altman, Large-scale labeling and assessment of sex bias in publicly available expression data, *BMC bioinformatics* **22**(1) (2021), 1–23.
- [57] R. Awadallah, M. Ramanath and G. Weikum, OpinioNetIt: A structured and faceted knowledge-base of opinions, in: *2012 IEEE 12th International Conference on Data Mining Workshops*, IEEE, 2012, pp. 878–881.
- [58] L. Gregori, R. Varvara and A.A. Ravelli, Action Type induction from multilingual lexical features, *Procesamiento del Lenguaje Natural* **63** (2019), 85–92.
- [59] A. Olteanu, C. Castillo, F. Diaz and E. Kıcıman, Social data: Biases, methodological pitfalls, and ethical boundaries, *Frontiers in Big Data* **2** (2019), 13.
- [60] C.A. Bail, L.P. Argyle, T.W. Brown, J.P. Bumpus, H. Chen, M.F. Hunzaker, J. Lee, M. Mann, F. Merhout and A. Volfovsky, Exposure to opposing views on social media can increase political polarization, *Proceedings of the National Academy of Sciences* **115**(37) (2018), 9216–9221.
- [61] C.T. Nguyen, Echo chambers and epistemic bubbles, *Episteme* **17**(2) (2020), 141–161.
- [62] F. Kamiran and T. Calders, Data preprocessing techniques for classification without discrimination, *Knowledge and information systems* **33**(1) (2012), 1–33.
- [63] S. Romero and K. Becker, Improving the classification of events in tweets using semantic enrichment, in: *Proceedings of the International Conference on Web Intelligence*, 2017, pp. 581–588.
- [64] A. Schulz, C. Guckelsberger and F. Janssen, Semantic abstraction for generalization of tweet classification: An evaluation of incident-related tweets, *Semantic Web* **8**(3) (2017), 353–372.
- [65] M. Wick, S. Panda and J.-B. Tristan, Unlocking fairness: a trade-off revisited (2019).
- [66] E.M. Bender, T. Gebru, A. McMillan-Major and S. Shmitchell, On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?, in: *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 2021, pp. 610–623.
- [67] S.A. Friedler, C. Scheidegger and S. Venkatasubramanian, On the (im) possibility of fairness, *arXiv preprint arXiv:1609.07236* (2016).

- [68] M.v. Erp and V.d. Boer, A Polyvocal and Contextualised Semantic Web, in: *European Semantic Web Conference*, Springer, 2021, pp. 506–512.
- [69] T.G. Papaioannou, V. Stankovski, P. Kochovski, A. Simonet-Boulogne, C. Barelle, A. Ciaramella, M. Ciaramella and G.D. Stamoulis, A New Blockchain Ecosystem for Trusted, Traceable and Transparent Ontological Knowledge Management, in: *International Conference on the Economics of Grids, Clouds, Systems, and Services*, Springer, 2021, pp. 93–105.
- [70] D. Stepanova, M.H. Gad-Elrab and V.T. Ho, Rule induction and reasoning over knowledge graphs, in: *Reasoning Web International Summer School*, Springer, 2018, pp. 142–172.
- [71] S. Mishra, S. He and L. Belli, Assessing Demographic Bias in Named Entity Recognition, *arXiv preprint arXiv:2008.03415* (2020).
- [72] A. Gaut, T. Sun, S. Tang, Y. Huang, J. Qian, M. ElSherief, J. Zhao, D. Mirza, E. Belding, K.-W. Chang et al., Towards understanding gender bias in relation extraction, *arXiv preprint arXiv:1911.03642* (2019).
- [73] J. Fisher, D. Palfrey, C. Christodoulopoulos and A. Mittal, Measuring social bias in knowledge graph embeddings, *arXiv preprint arXiv:1912.02761* (2019).
- [74] W. Radstok, M. Chekol, M. Schaefer et al., Are knowledge graph embedding models biased, or is it the data that they are trained on?, in: *Wikidata Workshop 2021 co-located with the 20th International Semantic Web Conference (ISWC 2021)*, 2021.
- [75] S. Tonelli, C. Giuliano and K. Tymoshenko, Wikipedia-based WSD for multilingual frame annotation, *Artificial Intelligence* **194** (2013), 203–221.
- [76] J. Teixeira, L. Sarmento and E. Oliveira, Comparing verb synonym resources for portuguese, in: *International Conference on Computational Processing of the Portuguese Language*, Springer, 2010, pp. 100–109.
- [77] A. Soulet, A. Giacometti, B. Markhoff and F.M. Suchanek, Representativeness of knowledge bases with the generalized Benford’s law, in: *International Semantic Web Conference*, Springer, 2018, pp. 374–390.
- [78] A.G. Nuzzolese, A. Gangemi, V. Presutti and P. Ciancarini, Type inference through the analysis of wikipedia links, in: *LDOW*, 2012.
- [79] N. Chah and P. Andritsos, WikiMetaData Studio: Dashboards From Data Profiling the Languages, Properties, and Items of Wikidata.
- [80] P. Vougiouklis, E. Maddalena, J. Hare and E. Simperl, How biased is your NLG evaluation? (2018).
- [81] A. Rossi, D. Barbosa, D. Firmani, A. Matinata and P. Meraldo, Knowledge graph embedding for link prediction: A comparative analysis, *ACM Transactions on Knowledge Discovery from Data (TKDD)* **15**(2) (2021), 1–49.
- [82] F. Akrami, L. Guo, W. Hu and C. Li, Re-evaluating embedding-based knowledge graph completion methods, in: *Proceedings of the 27th ACM international conference on information and knowledge management*, 2018, pp. 1779–1782.
- [83] B. Kruit, P. Boncz and J. Urbani, Extracting novel facts from tables for knowledge graph completion, in: *International semantic web conference*, Springer, 2019, pp. 364–381.
- [84] C.T. Wolf, From Knowledge Graphs to Knowledge Practices: On the Need for Transparency and Explainability in Enterprise Knowledge Graph Applications (2020).
- [85] C.M. Keet, An exploration into cognitive bias in ontologies (2021).
- [86] M.K. Halawani, R. Forsyth and P. Lord, A literature based approach to define the scope of biomedical ontologies: A case study on a rehabilitation therapy ontology, *arXiv preprint arXiv:1709.09450* (2017).
- [87] R. Mac An Tsaoir, Using Spreading Activation to Evaluate and Improve Ontologies, in: *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, 2014, pp. 2237–2248.
- [88] E.J. Gerritse, F. Hasibi and A.P. de Vries, Bias in conversational search: The double-edged sword of the personalized knowledge graph, in: *Proceedings of the 2020 ACM SIGIR on International Conference on Theory of Information Retrieval*, 2020, pp. 133–136.
- [89] A. Huizing, C. Veenman, M. Neerinx and J. Dijk, Hybrid AI: The Way Forward in AI by Developing Four Dimensions, in: *International Workshop on the Foundations of Trustworthy AI Integrating Learning, Optimization and Reasoning*, Springer, 2020, pp. 71–76.
- [90] .
- [91] A. Miles and S. Bechhofer, SKOS simple knowledge organization system reference, *W3C recommendation* (2009).
- [92] C. Manning and H. Schütze, *Foundations of statistical natural language processing*, MIT press, 1999.
- [93] F.A. Satti, M. Hussain, J. Hussain, S.I. Ali, T. Ali, H.S.M. Bilal, T. Chung and S. Lee, Unsupervised Semantic Mapping for Healthcare Data Storage Schema, *IEEE Access* **9** (2021), 107267–107278.
- [94] X. Wang, X. Han, Z. Liu, M. Sun and P. Li, Adversarial training for weakly supervised event detection, in: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 2019, pp. 998–1008.
- [95] Y. Xiang, Y. Zhang, X. Wang, Y. Qin and W. Han, Bias modeling for distantly supervised relation extraction, *Mathematical Problems in Engineering* **2015** (2015).