# Bias in Conversational Search:
# The Double-Edged Sword of the Personalized Knowledge Graph

Emma J. Gerritse
Radboud University
emma.gerritse@ru.nl

Faegheh Hasibi
Radboud University
f.hasibi@cs.ru.nl

Arjen P. de Vries
Radboud University
a.devries@cs.ru.nl

## ABSTRACT

Conversational AI systems are being used in personal devices, providing users with highly personalized content. Personalized knowledge graphs (PKGs) are one of the recently proposed methods to store users' information in a structured form and tailor answers to their liking. Personalization, however, is prone to amplifying bias and contributing to the echo-chamber phenomenon. In this paper, we discuss different types of biases in conversational search systems, with the emphasis on the biases that are related to PKGs. We review existing definitions of bias in the literature: people bias, algorithm bias, and a combination of the two, and further propose different strategies for tackling these biases for conversational search systems. Finally, we discuss methods for measuring bias and evaluating user satisfaction.

## CCS CONCEPTS

• **Information systems → Information retrieval**; **Users and interactive retrieval**; **Personalization**;

## KEYWORDS

Bias, Conversational search, Knowledge graphs, Fairness

## 1 INTRODUCTION

Conversational search has gained great attention recently, driven by the success of conversational AI systems such as Alexa, Google Home, and their likes. These systems are being installed on personal devices such as watches, phones, and laptops, collecting more and more personal information from users. Ideas have been proposed to make use of this personal information and build personalized systems, in order to gain improved user satisfaction. One of the proposed approaches is building *personalized knowledge graphs (PKGs)*,

which is a structured form of information about "entities personally related to a user, their attributes and the relations between them" [2]. PKGs provide the system with (locally stored) rich information about the user, to tailor the answers as much to their liking as possible [2]. However, introducing such a data resource and personalizing search results is not without risks: personalization can lead to echo-chamber effect and reinforce bias [4].

Bias in search and recommendation systems can influence human decision making, contribute to societal and political biases, and impact the health of our society [14]. Examples of bias in current search and social media websites include Google's anti-Brexit bias[1], Facebook's filter bubbles during US presidential election[2], and Linkedin's gender bias[3]. Bias can get even more severe in a conversational setting because: (i) answers in conversational systems are concise and there is less chance of showing diversified results compared to the traditional ten blue links interface, and (ii) users tend to access information that are in line with their prior views [12], and conversational systems provide personalized results to optimize for high user satisfaction and engagement, thereby intensifying bias. This raises several questions about bias in conversational search. Can forms of bias be introduced in this personalized search setting, which would not happen so easily in other search settings? How will conversational systems influence the behavior of their users? And, can users influence the behavior of their conversational agents? What are the sources of bias and how can we measure bias? How to mitigate bias while keeping users engaged and satisfied?

In this position paper, we discuss these questions from multiple angles. We consider the risk of different types of biases, and the places in the system architecture where these might be introduced; a major challenge being the difficulty to distinguish between *bias* and *preference*. We discuss how bias can be introduced by search algorithms, specifically when constructing PKGs from users' personal data and using them to provide personalized responses. We further explain how PKGs play the role of a double-edged sword to on one hand amplify bias but on the other hand detect and mitigate bias (Section 2).

Next, we argue that the community has to consider conversational search systems in regard to questions of bias, and find solutions to address it (Section 3). The main usage of a conversational system is a convenient form of tackling search; not an assistant to educate the user. If we would correct system results for bias too strongly, the debiasing process can only degrade user satisfaction.

---

[1] https://www.dailymail.co.uk/news/article-7605265/Google-facing-claims-anti-Brexit-bias-web-searches.html

[2] https://www.theguardian.com/us-news/2016/nov/16/facebook-bias-bubble-us-election-conservative-liberal-news-feed

[3] https://time.com/4484530/linkedin-gender-bias-search/

Conversational search therefore faces quite a challenge with respect to design ethics: not only do we need to think through the biases that we encounter and find solutions to counter those biases, we also have to do so in a manner that maintains user satisfaction. We outline possible approaches to help avoid the reinforcement of biases, while keeping the attraction of personalization to improve user satisfaction. Finally, we discuss methods to measure bias and user satisfaction (with respect to bias measures) in conversational systems (Section 4).

## 2 TYPES OF BIAS IN CONVERSATIONAL SEARCH

Baeza-Yates [1] enumerates different types of bias on the web and groups them into three categories: (i) bias that involves only algorithms, (ii) bias that originates from people, and (iii) bias that involves both algorithm and people. We discuss the different types of biases we expect to encounter in conversational search along the same categories.

### 2.1 Bias from Algorithms

*Knowledge graph construction bias.* Knowledge graphs (KGs) are structured repositories of data and powerful means to provide a machine understandable form of knowledge. It is envisaged that public, domain-specific, and personalized knowledge graphs (PKGs) can be used to empower conversational search systems [2]. While generation of general purpose knowledge graphs (e.g., YAGO and WikiData) and domain-specific knowledge graphs (e.g., GeoNames and MusicBrainz) can involve bias, construction of personalized knowledge graphs can be even a greater source of bias.

PKGs can be built from social media feeds, search history [11], conversation history, and other sources of information for which users give the system access permission (e.g., online shops and contacts). The choice of data source by itself can already introduce bias in PKGs. If the system relies on publicly available social media feeds (e.g., Twitter [18]), PKGs may be biased towards only one aspect of users' preferences: people often use different social media websites for different purposes; e.g., Twitter for work and Instagram for casual usage. Even if the system gets access to the content of a user's private social media feeds, the question remains whether the PKG is diversified enough or not: the user might avoid sharing some aspects of her life in social media. Therefore, conversational systems need to account for the incompleteness of PKGs and incorporate it into their search and recommendation algorithms. Otherwise, this can negatively affect users' satisfaction, as the user encounter undesired bias in search and recommendation results (e.g., observing only work-related events instead of diversified ones).

Another source of bias in constructing PKGs is *time* [2]. Time-based events can greatly influence search behavior. For example, during election times people are more likely to share their opinion about politics, even when normally politics is not an important aspect of their lives. Another example is that during the Covid-19 outbreak in 2020, many people used social media to discuss the epidemiology, even though epidemiology is normally not in their interests. When using social media for knowledge graph construction, this can lead to some kind of time-based bias or time-based clutter, which will "somehow" have to be cleaned up after the event.

*Search algorithms bias.* Bias can also occur in search and recommendation algorithms. While not much research covers the biases in conversational search, some research has been carried out to measure bias in neural ranking models and embedding algorithms. In a recent study, Rekabsaz and Schedl [16] found that gender bias in document ranking is intensified when using neural models and in particular the ones that are based on contextual embeddings like BERT [8].

Bias has been also observed in the word and graph embeddings themselves; e.g., social bias in WikiData graph embeddings [9] and gender bias in word embeddings (like Word2Vec) [3]. While debiasing embeddings may appear as an immediate solution to mitigate bias, incorporating debiased embeddings in information retrieval algorithms may not immediately affect the ranking outcome, as confirmed by some initial results in [10].

### 2.2 Bias from People

People with their activities and preferences themselves form another source of bias in conversational systems. A study on age demographic of conversational systems' users in the United States showed that the most active users are 30-44 years old [4]. Indeed, people's preferences vary by their age group, and therefore a large population of a certain age group can introduce certain biases in the system. For example, a study on recommender systems [19] shows that that preference of phone color differs per age group, and recommendations can be tailored to the age group of users. Not only age group, but also gender, location, and language of people can affect the outcome of conversational systems, often causing results biased towards the majority groups.

### 2.3 Bias from Algorithms and People

Conversational systems are highly interactive and are often designed to be personalized. These two features make conversational systems prone to a vicious cycle of bias that is generated by a biased user interacting with a biased algorithm, thereby creating a filter-bubble effect. Personalized knowledge bases are double-edged swords in such a setting: they can be used to accelerate personalization and therefore providing even more biased results to the user, or they can be used to detect a biased user and help to diversify results. Drawing a line between personalization and diversification is a challenge and varies for different scenarios. Here, we enumerate three different scenarios as examples:

- A user wants to order the same guitar strings (s)he previously ordered. However, these might not be the highest quality strings. Should the system recommend other strings or guide the user to buy the previously bought strings?
- A user has a certain political preference and asks a political question about her favorite party. Should the system takes users' preferences into account, or provide a neutral answer?
- A user searches for a known conspiracy, e.g., fake news about Covid-19 outbreak. Knowing that the user trusts the source of information based on available information in PKG, how the system should respond to the query?

In the following section, we propose some solutions that a system can take to handle bias.

---

[4]https://voicebot.ai/2019/06/21/voice-assistant-demographic-data-young-consumers-more-likely-to-own-smart-speakers-while-over-60-bias-toward-alexa-and-siri/

## 3 POSSIBLE STRATEGIES

Conversational search systems can take a wide range of strategies against bias, from being completely ignorant to chastising users. In this section, we discuss three main strategies that these systems can take and argue against being ignorant.

### 3.1 Ignorant

In this setting, the system does not detect bias and takes no action against it. Depending on the search algorithm, the system may provide nonsensical responses or get along with users and provide them with what they ask. Such systems can fall into the trap of creating a vicious cycle of second-order bias (generated from people and algorithmic bias) [1] and then become their own enemy. An example of such a scenario is Microsoft's AI chatbot that had to be taken down 16 hours after release, because of learning harmful intents from some Twitter users[5]. We argue that being ignorant and *not accounting for bias* is not an option for conversational search systems, especially because they are meant to work in a highly personalized setting.

### 3.2 User in Control

In this strategy, systems detect bias and offer the user control through its UI/UX, from being viciously biased or offering diversified results; see Figure 1 for an example. There is a wide spectrum of actions that systems take in this strategy, ranging from complete neutralization of results to playing along with users until reaching a certain threshold (e.g., inappropriate user requests). Choosing a spot in this wide spectrum of approaches is a challenge and requires carrying out user studies. In addition, users may have different tolerance for interacting with diversified answers: some may appreciate it and others may abandon the system altogether. Most of the current commercial systems choose to "play along" with users or take a middle approach of not answering or providing "Do not know" answers to gender and sexual questions [5]. This, however, may not be an ideal setting for the well being of all users. Providing users with an option to choose a degree of diversification (and therefore mitigate bias), e.g. with a UI control like a slider, may alleviate these issues. This does however raise additional implementation challenges, especially for data-driven approaches (like neural networks). Incorporating knowledge graphs and in particular PKGs can help detecting bias in the first place and further control it by providing factual and diversified responses (e.g., alternating between nodes of PKGs) to reduce bias.

### 3.3 Reflective

Another approach is being reflective about bias: both user and system detect bias and explicitly reflect it to each other. So here debiasing is the responsibility of both user and system.

In the case of user bias, the system notifies the user in an appropriate language and suggests alternatives (see Figure 2). Educating users with an explicit response may be considered as an over-reaction, however, it is desired for some cases; e.g., health-related fake news in a pandemic like COVID-19, sexual harassment,
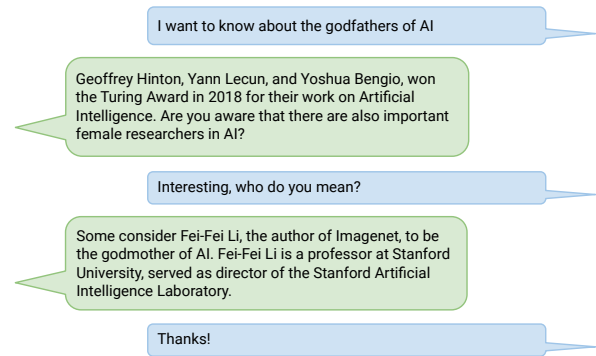
---

[5]https://www.bbc.com/news/technology-35890188



**Figure 1: Example of the *User in Control* strategy. The system detects the possible bias in the word 'godfathers', and offers the user also some interesting female answers.**
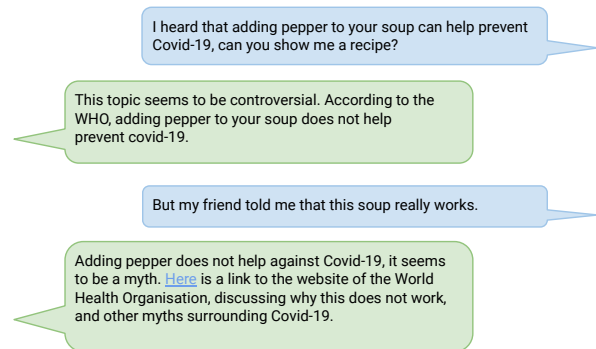


**Figure 2: Example of the *Reflective* strategy, where the system detects that the user is biased towards fake news, and the system supplies them with an alternative.**

or cases when the user is gathering confrontational viewpoints for an argument.

Users may also correct an unnecessary bias of system and help to create a balance between bias and preference. Examples of such scenarios include: users changing their interests and opinions over time, or searching for something either on the spur of the moment or on someone else's request. In such cases, systems should provide the user with the source of bias and the option to mediate it. PKGs, while being a source of bias, can help the explainability of the system and enable users to influence the behavior of systems by editing their own PKG.

## 4 EVALUATION

In this section, we propose methods to measure bias in conversational search and further evaluate user satisfaction when applying

debiasing methods. Here, we do not aim to provide formal evaluation measures for bias, focusing instead on methods that can identify bias.

## 4.1 Measuring bias

Researchers have raised challenges of measuring bias in online information [15] and search engines [13]. Measuring bias in conversational search raises even more challenges, due to the difficulties of evaluating the systems themselves. While test collections like TREC CAST [7] and QuAC [6] exist, they are often restricted to the question-answering part of conversational systems. Evaluating multiple aspects of a conversational system is often performed with user experiments, and is mainly performed by considering general user satisfaction. We argue that a system's ability to handle bias should be measured separately and be an integral element of user evaluations in personalized conversational search systems.

Another option to measure bias in conversational systems is to have an evaluation method that does not require relevance judgments. In [16], one such method is discussed for computing gender bias in search results. This is performed by computing a term frequency based score with predefined gendered words. Similar approaches can be employed to measure the bias of an answer to a query, not only with respect to gender but other aspects like political bias.

It is important to note that bias is a relative concept and depending on the context, a biased answer may not be always undesired. One way of measuring bias is computing the results with and without PKGs; e.g., by using the general purpose knowledge graphs. If the difference is rather large, then it indicates that personalization has occurred in a large degree (which is not always an undesired behavior). This is indeed under the assumption that global knowledge graphs are fair representation of information. One can also generate the results with an *opposite knowledge graph*, where nodes and relations of user's PKG are replaced with opposite information. We acknowledge that this concept does not generalize, as opposite information may not be available for all nodes. However, this solution may provide interesting insights for certain biases like gender and politics. Consider for example swapping "Female", "right-wing", and "conservative" with "Male", "left-wing", and "progressive", and compare the bias scores as described in [16]. This indicates the influence of these nodes on the results. If there is a large difference between the two polarities, it means that personalization has occurred to a large extent.

## 4.2 Measuring user-satisfaction

Changing system behavior to handle bias, either by diversifying results or offering a reflective option, may change user satisfaction. It is expected that users do not appreciate the feeling of being judged or educated by their conversational assistant. It is therefore wise to first study the best strategy of handling bias, before investigating technical aspects of offering these solutions. To this end, users satisfaction can be measured by wizard experiments, similar to [17]. Here, a wizard is a human who pretends to be a conversational search system and interacts with a human. These studies can be performed by instructing a wizard to offer perspectives to questions of human test subjects and measure user satisfaction.

## 5 CONCLUSION

In this paper, we touched upon some aspects of bias in conversational search. We raised questions about bias in the personalized search setting, especially using personalized knowledge graphs. We discussed the challenge of distinguishing between bias and preference, caused by the inherent bias of user preferences and also limitations of offering diversified results in conversational settings. We enumerated different types of bias present in conversational search, being algorithmic bias, people bias and bias introduced by both algorithm and people. We also suggested possible solutions for handling bias, based on diversification and mutual reflection of bias between user and system. Lastly, we suggested methods of measuring bias and evaluating user satisfaction. Bias in conversational search is an aspect that requires extra attention and we urge the community to keep it in mind while developing conversational search systems.

## REFERENCES

[1] Ricardo Baeza-Yates. 2018. Bias on the Web. *Commun. ACM* 61, 6 (2018), 54–61.
[2] Krisztian Balog and Tom Kenter. 2019. Personal Knowledge Graphs: A Research Agenda. In *Proc. of ICTIR '19*. 217–220.
[3] Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. 2016. Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings. In *Proc. of NIPS '16*. 4356–4364.
[4] Engin Bozdag. 2013. Bias in Algorithmic Filtering and Personalization. *Ethics and Inf. Technol.* 15, 3 (2013), 209–227.
[5] Amanda Cercas Curry and Verena Rieser. 2018. #MeToo Alexa: How Conversational Systems Respond to Sexual Harassment. In *Proc. of the Second ACL Workshop on Ethics in Natural Language Processing*. 7–14.
[6] Eunsol Choi, He He, Mohit Iyyer, Mark Yatskar, Wen-tau Yih, Yejin Choi, Percy Liang, and Luke Zettlemoyer. 2018. QuAC: Question Answering in Context. In *Proc. of EMNLP '18*. 2174–2184.
[7] Jeffrey Dalton, Chenyan Xiong, and Jamie Callan. 2020. TREC CAsT 2019: The Conversational Assistance Track Overview. *CoRR* abs/2003.13624 (2020).
[8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proc. of NAACL '19*. 4171–4186.
[9] Joseph Fisher. 2019. Measuring Social Bias in Knowledge Graph Embeddings. *arXiv preprint arXiv:1912.02761* (2019).
[10] Emma Gerritse and Arjen de Vries. 2020. Effect of Debiasing on Information Retrieval. In *Proc. of the International Workshop on Algorithmic Bias in Search and Recommendation*.
[11] Jiyin He and Marc Bron. 2017. Measuring Demonstrated Potential Domain Knowledge with Knowledge Graphs. In *Proc. of the First Workshop on Knowledge Graphs and Semantics for Text Retrieval and Analysis*, Vol. 1883. 13–18.
[12] Danai Koutra, Paul N Bennett, and Eric Horvitz. 2015. Events and Controversies: Influences of a Shocking News Event on Information Seeking. In *Proc. of WWW '15*. 614–624.
[13] Abbe Mowshowitz and Akira Kawaguchi. 2005. Measuring Search Engine Bias. *Inf. Process. Manage.* 41, 5 (2005), 1193–1205.
[14] Alexandra Olteanu, Jean Garcia-Gathright, Maarten de Rijke, and Michael D Ekstrand. 2019. FACTS-IR: Fairness, Accountability, Confidentiality, Transparency, and Safety in Information Retrieval. *SIGIR Forum* 53, 2 (2019), 20–43.
[15] Evaggelia Pitoura, Panayiotis Tsaparas, Giorgos Flouris, Irini Fundulaki, Panagiotis Papadakos, Serge Abiteboul, and Gerhard Weikum. 2018. On Measuring Bias in Online Information. *SIGMOD Rec.* 46, 4 (2018), 16–21.
[16] Navid Rekabsaz and Markus Schedl. 2020. Do Neural Ranking Models Intensify Gender Bias?. In *Proc. of SIGIR '20*.
[17] Alexandra Vtyurina, Denis Savenkov, Eugene Agichtein, and Charles L. A. Clarke. 2017. Exploring Conversational Search With Humans, Assistants, and Wizards. In *Proc. of CHI '17*. 2187–2193.
[18] An-Zi Yen, Hen-Hsen Huang, and Hsin-Hsi Chen. 2019. Personal Knowledge Base Construction from Text-based Lifelogs. In *Proc. of SIGIR '19*. 185–194.
[19] Xin Wayne Zhao, Yanwei Guo, Yulan He, Han Jiang, Yuexin Wu, and Xiaoming Li. 2014. We Know What You Want to Buy: A Demographic-based System for Product Recommendation on Microblogs. In *Proc. of SIGKDD '14*. 1935–1944.