

# Evaluation Metrics for Measuring Bias in Search Engine Results

Gizem Gezici<sup>1</sup>, Aldo Lipani<sup>2</sup>, Yucel Saygin<sup>1</sup>,  
and Emine Yilmaz<sup>2</sup>

<sup>1</sup>Sabanci University, Department of Computer  
Science and Engineering, Istanbul, Turkey

<sup>2</sup>University College London, Department of  
Computer Science, London, UK

Received: 4 February 2020 / Accepted: 7 December 2020

**Abstract** Search engines decide what we see for a given search query. Since many people are exposed to information through search engines, it is fair to expect that search engines are neutral. However, search engine results do not necessarily cover all the viewpoints of a search query topic, and they can be biased towards a specific view since search engine results are returned based on relevance, which is calculated using many features and sophisticated algorithms where search neutrality is not necessarily the focal point. Therefore, it is important to evaluate the search engine results with respect to bias. In this work we propose novel web search bias evaluation measures which take into account the rank and relevance. We also propose a framework to evaluate web search bias using the proposed measures and test our framework on two popular search engines based on 57 controversial query topics such as abortion, medical marijuana, and gay marriage. We measure the *stance bias* (in support or against), as well as the *ideological bias* (conservative or liberal). We observe that the stance does not necessarily correlate with the ideological leaning, e.g. a positive stance on abortion indicates a liberal leaning but a positive stance on Cuba embargo indicates a conservative leaning. Our experiments show that neither of the search engines suffers from stance bias. However, both search engines suffer from ideological bias, both favouring one ideological leaning to the other, which is more significant from the perspective of polarisation in our society.

**Keywords** Bias evaluation, Fair ranking, Search bias, Web Search

---

✉ Gizem Gezici  
gizemgezici@sabanciuniv.edu

Address(es) of author(s) should be given

## 1 Introduction

Search engines have become an indispensable part of our lives. As reported by SmartSights (2018), 46.8% of the world population accessed the internet in 2017 and by 2021, the number is expected to reach 53.7%. According to InternetLiveStats (2018), currently on average 3.5 billion Google searches are done per day. These statistics indicate that search engines replaced traditional broadcast media and have become a *major* source of information “gatekeepers to the Web” for many people (Diaz, 2008). As information seekers search the Web more, they are also more influenced by Search Engine Result Pages (SERPs), pertaining to a wide range of areas (e.g., work, entertainment, religion, and politics). For instance, in the course of elections, it is known that people issue repeated queries on the Web about political candidates and events such as “democratic debate”, “Donald Trump”, “climate change” (Kulshrestha *et al.*, 2018). SERPs returned in response to these queries may influence the voting decisions as claimed by Epstein & Robertson (2015), who report that *manipulated* search rankings can change the voting preferences of undecided individuals at least by 20%.

Although search engines are widely used for seeking information, the majority of online users tend to believe that they provide *neutral* results, i.e. serving only as facilitators in accessing information on the Web (Goldman, 2008). However, there are counter examples to that belief as well. A recent dispute between the U.S. President Donald Trump and Google is such an example, where Mr. Trump accused Google of displaying only negative news about him when his name is searched to which Google responded by saying: “When users type queries into the Google Search bar, our goal is to make sure they receive the most relevant answers in a matter of seconds” and “Search is not used to set a political agenda and we don’t bias our results toward any political ideology” (Ginger & David, 2018). In this work, we hope to shed some light on that debate, by not specifically concentrating on queries regarding Donald Trump but by conducting an in depth analysis of search answers to a broad set of controversial topics based on concrete evaluation measures.

*Bias* is defined with respect to balance in representativeness of Web documents retrieved from a database for a given query (Mowshowitz & Kawaguchi, 2002a). When a user issues a query to a search engine, documents from different sources are gathered, ranked, and displayed to the user. Assume that a user searches for *2016 presidential election* and the top-n ranked results are displayed. In such a search scenario, the retrieved results may favor some political perspectives over others and thereby fail to provide impartial knowledge for the given query as claimed by Mr. Trump, though without any scientific support. Hence, the potential *undue emphasis* of specific perspectives (or viewpoints) in the retrieved results lead to bias (Kulshrestha *et al.*, 2018). With respect to the definition of bias and the presented scenario, if there is an unbalanced representation, i.e. skewed or slanted distribution, of the viewpoints in a SERP, i.e. not only in political searches, towards the query’s topic, then we consider this SERP as *biased* for the given search query.

77 Bias is especially important if the query topic is *controversial* having op-  
78 posing views, in which case it becomes more critical that search engines are  
79 supposed to return results with a *balanced* representation of different perspec-  
80 tives which implies that they do not favour one specific perspective over an-  
81 other. Otherwise, this may dramatically affect public as in the case of elections  
82 leading to polarisation in society for *controversial* issues. On the other hand,  
83 returning an unbalanced representation of distinct viewpoints is not sufficient  
84 to claim that the search engine’s ranking algorithm is biased. One reason for a  
85 skewed SERP could be due to the corpus itself, i.e. if documents indexed and  
86 returned for a given topic come from a slanted distribution, meaning that the  
87 ranking algorithm returns a biased result set due to a biased corpus. To differ-  
88 entiate the algorithmic vs corpus bias, one needs to investigate the source of  
89 bias in addition to the skewed list analysis of the top-n search results. However,  
90 the existence of bias, regardless of being corpus or algorithmic bias, would still  
91 conflict with the expectation that an IR system should be fair, accountable,  
92 and transparent (Culpepper *et al.*, 2018). Furthermore, it was reported that  
93 people are more susceptible to bias when they are unaware of it (Bargh *et al.*,  
94 2001), and Epstein *et al.* (2017) showed that alerting users about bias can  
95 be effective in suppressing search engine manipulation effect (SEME). Thus,  
96 search engines should at least inform their users about the bias and decrease  
97 the possible SEME by making themselves more accountable, thereby alleviat-  
98 ing the negative effects of bias and serving only as facilitators as they generally  
99 claim to be. In this work, we aim to serve that purpose by proposing a search  
100 bias evaluation framework taking into account the rank and relevance <sup>1</sup> of the  
101 SERPs. Our contributions in this work can be summarised as follows:

- 102 1. We propose a *new generalisable search bias evaluation framework* to mea-  
103 sure bias in SERPs by quantifying two different types of bias on content  
104 which are stance bias and ideological bias.
- 105 2. We present *three novel fairness-aware measures of bias* that do not suffer  
106 from the limitations of the previously presented bias measures, based  
107 on common Information Retrieval (IR) *utility-based* evaluation measures:  
108 Precision at cut-off (P@n), Rank Biased Precision (RBP), and Discounted  
109 Cumulative Gain at cut-off (DCG@n) which are explained in Section 3.2  
110 in detail.
- 111 3. We apply the proposed framework to *measure the stance and ideological*  
112 *bias* not only in political searches but searches related to a wide range of  
113 controversial topics; including but not limited to education, health, enter-  
114 tainment, religion and politics on Google and Bing *news* search results.
- 115 4. We also utilise our framework to *compare the relative bias* for queries from  
116 various controversial issues on two popular search engines: Google and Bing  
117 news search.

118 We would like to note that we distinguish the stance and ideological leaning  
119 in SERPs. The stance in a SERP for a query topic could be in favor or against

---

<sup>1</sup> We are referring to the notion of relevance defined in the literature as system relevance, or topical relevance which is the relevance predicted by the system.

120 the topic, whereas the ideological leaning in a SERP stands for the specific  
121 ideological group as conservatives or liberals that supports the corresponding  
122 topic. Hence, the stance in a SERP does not directly imply the ideological  
123 leaning. For example, given two controversial queries, "abortion" and "Cuba  
124 embargo", a SERP could have a positive stance for the topic of abortion, indi-  
125 cating a liberal leaning, while a positive stance for the topic of Cuba embargo  
126 indicates a conservative leaning. Therefore looking at the stance of the SERPs  
127 for controversial issues is not enough and could even be misleading in deter-  
128 mining the ideological bias. We demonstrate how the proposed framework can  
129 be used to quantify bias in the SERPs of search engines (in this case Bing and  
130 Google) in response to queries related to *controversial* topics. Our analysis is  
131 mainly two-fold where we first evaluate stance bias in SERPs, and then use  
132 this evaluation as a proxy to quantify ideological bias asserted in the SERPs  
133 of the search engines.

134 In this work, via the proposed framework, we aim to answer the following  
135 research questions:

136 RQ1: On a pro-against stance space, do search engines return *biased* SERPs  
137 towards controversial topics?

138 RQ2: Do search engines show *significantly different* magnitude of stance bias  
139 from each other towards controversial topics?

140 RQ3: On a conservative-liberal ideology space, do search engines return *biased*  
141 SERPs and if so; are these biases *significantly different* from each other  
142 towards controversial topics?

143 We address these research questions for controversial topics representing a  
144 broad range of issues in SERPs of Google and Bing through content analy-  
145 sis, i.e. analysing the textual content of the retrieved documents. In order to  
146 answer RQ1, we measure the degree of deviation of the ranked SERPs from  
147 an *ideal* distribution, where different stances are *equally* likely to appear. To  
148 detect bias which results from the unbalanced representation of distinct per-  
149 spectives, we label the documents' stances with crowd-sourcing and use these  
150 labels for stance bias evaluation. In this paper we focus on a particular kind  
151 of bias, *statistical parity* or more generally known as *equality of outcome*, i.e.  
152 given a population divided into groups, the groups in the output of the sys-  
153 tem should be equally represented. This is in contrast with the other popular  
154 measure generally known as *equality of opportunity*, i.e. given a population  
155 divided into groups, the groups in the output should be represented based  
156 on their proportion in the population namely, base rates. For choosing the  
157 *equality of outcome*, we have mainly two reasons. First, in the context of the  
158 controversial topics, not all of the corresponding debate questions (queries)  
159 have certain answers based on scientific facts. Second, the identification of the  
160 stance for the full ranking list, i.e. which is a fair representative set of the in-  
161 dexed documents, is too expensive to get annotated through crowd-sourcing.  
162 Thus, this choice of *ideal* ranking makes the experiments feasible. To address  
163 RQ2, we compare the stance bias in the SERPs of the two search engines  
164 to see if they show similar level of bias for the corresponding controversial

165 topics. RQ3 is naturally answered by assigning an ideological leaning label to  
166 each query topic as conservative or liberal depending on which ideology favors  
167 the proposition in the query. We further interpret the document stance labels  
168 in conservative-to-liberal ideology <sup>2</sup> space and transform these stance labels  
169 into ideological leanings according to the assigned leaning labels of the corre-  
170 sponding topics. We note that conservative-to-liberal ideology space does not  
171 only stand for political parties. In this context, we accept these ideology labels  
172 as having a more conservative/liberal viewpoint towards a given controversial  
173 topic as similarly fulfilled by Lahoti *et al.* (2018) for three popular controversial  
174 topics of *gun control*, *abortion*, and *obamacare* in Twitter domain.

175 For instance, the topic of *abortion* has the query of *Should Abortion Be Le-*  
176 *gal?* Since mostly liberals support the proposition in this query, liberal leaning  
177 is assigned to abortion. The stance labels of the retrieved documents towards  
178 the query are transformed into ideological leanings as follows. If a document  
179 has the pro stance which means that it supports the asserted proposition,  
180 then its ideological leaning is liberal; if it has the against stance, its leaning is  
181 conservative.

182 In our bias evaluation framework, we concentrate on the top-10 SERPs  
183 coming from the *news* sources to investigate two major search engines (Bing  
184 and Google) in terms of bias. We deliberately use *news* SERPs for our ex-  
185 periments since they often exhibit a specific view towards a topic (Alam &  
186 Downey, 2014). Recent studies (Sarcona, 2019; 99Firms, 2019) show that on  
187 average more than 70% of all the clicks are in the first page results, thus we only  
188 focus on the top-10 results to show the existence of bias. Experiments show  
189 that there is no statistically significant difference of *stance bias* in magnitude  
190 measured across the two search engines, meaning that they do not favour one  
191 specific stance over other. However, we should stress that stance bias results  
192 need to be taken with a grain of salt as demonstrated through the abortion  
193 and Cuba embargo query examples. Polarisation of the society is mostly on  
194 ideological leanings, and our second phase of experiments show that there is  
195 statistically significant difference of *ideological bias*, where both search engines  
196 favour one ideological leaning over other.

197 The remainder of the paper is structured as follows. In Section 2 we give the  
198 related work and the search bias evaluation framework is proposed in Section  
199 3. In Section 4 we detail the experimental setup, and present the results. Then,  
200 we discuss the results in Section 5. In Section 6 we present the limitations of  
201 this work, and we conclude in Section 7.

## 202 2 Background & Related Work

203 In recent years, bias analysis in SERPs of search engines has attracted a lot of  
204 interest (Baeza-Yates, 2016; Mowshowitz & Kawaguchi, 2002b; Noble, 2018;  
205 Pan *et al.*, 2007; Tavani, 2012) due to the concerns that search engines may

---

<sup>2</sup> We are referring to the notion of ideology perceived by the crowd workers.

manipulate the search results influencing users. The main reason behind these concerns is that search engines have become the fundamental source of information (Dutton *et al.*, 2013), and surveys from Pew (2014) and Reuters (2018) found that more people obtain their news from search engines than social media. The users reported higher trust on search engines for the accuracy of information (Newman *et al.*, 2018, 2019; Elisa Shearer, 2018) and many internet-using US adults even use search engines to fact-check information (Dutton *et al.*, 2017).

To figure out how this growing usage of search engines and trust in them might have undesirable effects on public, and what could be the methods to measure those effects, in the following we review the research areas related first to automatic stance detection, then to fair ranking evaluation, and lastly to search bias quantification.

## 2.1 Opinion Mining and Sentiment Analysis

A form of Opinion Mining related to our work is Contrastive Opinion Modeling (COM). Proposed by Fang *et al.* (2012), in COM, given a political text collection, the task is to present the opinions of the distinct perspectives on a given query topic and to quantify their differences with an unsupervised topic model. COM is applied on debate records and headline news. Differently from keyword analysis to differentiate opinions using topic modelling, we compute different IR metrics from the content of the news articles to evaluate and compare the bias in the SERPs of two search engines. Aktolga & Allan (2013) consider the sentiment towards controversial topics and propose different diversification methods based on the topic sentiment. Their main aim is to diversify the retrieved results of a search engine according to various sentiment biases in blog posts rather than measure bias in the SERPs of *news* search engines as we do in this work.

Demartini & Siersdorfer (2010) exploit automatic and lexicon-based text classification approaches, Support Vector Machines and SentiWordNet respectively to extract sentiment value from the textual content of SERPs in response to controversial topics. Unlike us, Demartini & Siersdorfer (2010) only use this sentiment information to compare opinions in the retrieved results of three commercial search engines without measuring bias. In this paper, we propose a new bias evaluation framework with robust bias measures to systematically measure bias in SERPs. Chelaru *et al.* (2012) focus on queries rather than SERPs and investigate if the opinionated queries are issued to search engines by computing the sentiment of suggested queries for controversial topics. In a follow-up work (Chelaru *et al.*, 2013), authors use different classifiers to detect the sentiment expressed in queries and extend the previous experiments with two different use cases. Instead of queries, our work analyses the SERPs in *news* domain, therefore we need to identify the stance of the news articles. Automatically obtaining article stances is beyond the scope of this work, thus we use crowd-sourcing.

## 249 2.2 Evaluating Fairness in Ranking

250 Fairness evaluation in ranked results has attracted attention in recent years.  
 251 Yang & Stoyanovich (2017) propose three bias measures, namely Normalized  
 252 discounted difference (rND), Normalized discounted Kullback-Leibler diver-  
 253 gence (rKL) and Normalized discounted ratio (rRD) that are related to Nor-  
 254 malized Discounted Cumulative Gain (NDCG) through the use of logarithmic  
 255 discounting for regularization which is inspired from NDCG as also stated in  
 256 the original paper. Researchers use these metrics to check if there exists a sys-  
 257 tematic discrimination against a group of individuals, when there are only two  
 258 different groups as a protected ( $g_1$ ) and an unprotected group ( $g_2$ ) in a rank-  
 259 ing. In other words, researchers quantify the relative representation of  $g_1$  (the  
 260 protected group), whose members share a characteristic such as race or gender  
 261 that cannot be used for discrimination, in a ranked output. The definitions of  
 262 these three proposed measures can be rewritten as follows:

$$263 \quad f_{g_1}(r) = \frac{1}{Z} \sum_{i=10,20,\dots}^{|r|} \frac{1}{\log_2 i} |d_{g_1}(i, r)|, \quad (1)$$

264 where  $f(r)$  is a general definition of an evaluation measure for a given ranked  
 265 list of documents, i.e. a SERP, whereas  $f_{g_1}$  is specifically for the protected  
 266 group of  $g_1$ . In this definition,  $Z$  is a normalisation constant,  $r$  is the ranked  
 267 list of the retrieved SERP and  $|r|$  is the size of this ranked list, i.e. number of  
 268 documents in the ranked list. Note that,  $i$  is deliberately incremented by 10,  
 269 to compute *set-based fairness* at discrete values as top-10, top-20 etc., instead  
 270 of 1 as usually done in IR for the proposed measures to show the correct  
 271 behaviour with bigger sample sizes. The purpose of computing the *set-based*  
 272 *fairness* to express that being fair at higher positions of the ranked list is more  
 273 important, e.g. top-10 vs. top-100.

274 In the rewritten formula,  $d_{g_1}$  defines a distance function between the ex-  
 275 pected probability to retrieve a document belonging to  $g_1$ , i.e. in the overall  
 276 population, and its observed probability at rank  $i$  to measure the systematic  
 277 bias. These probabilities turn out to be equal to P@n:

$$278 \quad P_{g_1}@n = \frac{1}{n} \sum_{i=1}^n [j(r_i) = g_1], \quad (2)$$

279 when computed over  $g_1$  at cut-off value  $|r|$  and  $i$  for the three proposed mea-  
 280 sures as below. In this formula,  $n$  is the number of documents considered in  $r$   
 281 as a cut-off value, and  $r_i$  is defined as the document in  $r$  retrieved at rank  $i$ .  
 282 Note that,  $j(r_i)$  returns the label associated to the document  $r_i$  specifying its  
 283 group as  $g_1$  or  $g_2$ . Based on this,  $[j(r_i) = g_1]$  refers to a conditional statement  
 284 which returns 1 if the document  $r_i$  is the member of  $g_1$  and 0 otherwise. In  
 285 the original paper,  $d_{g_1}$  is defined for rND, rKL, and rRD as:

$$d_{g_1}(i, r) = P_{g_1 @ i} - P_{g_1 @ |r|} \quad \text{for rND,}$$

$$d_{g_1}(i, r) = -P_{g_1 @ i} \log \left( \frac{P_{g_1 @ |r|}}{P_{g_1 @ i}} \right) - (1 - P_{g_1 @ i}) \log \left( \frac{1 - P_{g_1 @ |r|}}{1 - P_{g_1 @ i}} \right) \quad \text{for rKL,}$$

$$d_{g_1}(i, r) = \frac{P_{g_1 @ i}}{1 - P_{g_1 @ i}} - \frac{P_{g_1 @ |r|}}{1 - P_{g_1 @ |r|}} \quad \text{for rRD.}$$

These measures, although inspired by IR evaluation measures, particularly in the context of content bias in search results suffer from the following limitations:

1. rND measure focuses on the protected group ( $g_1$ ). If we were to compute  $f$  at steps of 1 with the given equal desired proportion of the two groups as 50:50, then the distance function of rND, denoted as  $d_{g_1}$  would always give a value of 0.5 for the first retrieved document, where  $i = 1$ . This will always be the case, no matter which group this document belongs to, e.g. *pro* or *against* in our case. This is caused by  $d_{g_1}$  of rND through the use of its absolute value in Eq. (1). In our case, this holds when  $i = 1, 2, 4$  and  $r = 10$  where we measure bias in the top-10 results. This is in fact avoided in the original paper (Yang & Stoyanovich, 2017) by computing  $f$  at steps of 10 as top-10, top-20 etc. rather than the steps of 1 as it is usually done in IR which gives more meaningful results in our evaluation framework.
2. rKL measure cannot differentiate between biases of equal magnitude, but in opposite directions with the given equal desired proportion of the two groups as 50:50, i.e. it cannot differentiate bias towards *conservative*, or *liberal* in our case. Also, in IR settings it is not as easy to interpret the computed values from the KL-divergence (denoted as  $d_{g_1}$  for rKL) compared to our measures since our measures are based on the standard utility-based IR measures. Furthermore, KL-divergence tends to generate larger distances for small datasets, thus it could compute larger bias values in the case of only 10 documents, and this situation may become even more problematic if we measure bias for less number of documents, e.g. top-3, top-5 for a more fine-grained analysis. In the original paper, this disadvantage is alleviated by computing the rKL values also at discrete points of steps 10 instead of 1.
3. rRD measure does not treat the protected and unprotected groups ( $g_1$  and  $g_2$ ) symmetrically as stated in the original paper, which is not applicable to our framework. Our proposed measures treat  $g_1$  and  $g_2$  equal since we have two protected groups; *pro* and *against* for stance bias, *conservative* and *liberal* for ideological bias to measure bias in search settings. Moreover, rRD is only applicable in special conditions when  $g_1$  is the minority group



327 in the underlying population as also declared by the authors, while we  
328 do not have such constraints for our measures in the scope of search bias  
329 evaluation.

- 330 4. These measures focus on differences in the relative representation of  $g_1$  be-  
331 tween distributions. Therefore, from a general point of view, most probably  
332 more samples are necessary for these measures to show the expected behav-  
333 ior and work properly. In the original paper, experiments are fulfilled with  
334 three different datasets, one is synthetic which includes 1000 samples and  
335 two are real datasets which include 1000 and 7000 samples to evaluate bias  
336 with these measures, while we have only 10 samples for query-wise eval-  
337 uation. This is probably because these measures were mainly devised for  
338 the purpose of measuring bias in ranked outputs instead of search engine  
339 results; none of these datasets contain search results either.
- 340 5. These measures are difficult to use in practice, since they rely on a normal-  
341 ization term,  $Z$  that is computed stochastically, i.e. as the highest possible  
342 value of the corresponding bias measure for the given number of docu-  
343 ments  $n$  and protected group size  $|g_1|$ . In this paper, we rely on standard  
344 statistical tests, since they are easier to interpret, provide confidence inter-  
345 vals, and have been successfully used to investigate inequalities in search  
346 systems previously by Chen *et al.* (2018).
- 347 6. These measures do not consider relevance which is a fundamental aspect  
348 when evaluating bias in search engines. For example, as in our case, when  
349 searching for a controversial topic, if the first retrieved document is about  
350 a news belonging to  $g_1$  but its content is not relevant to the searched topic,  
351 then these measures would still consider this document as positive for  $g_1$ .  
352 However, this document has absolutely no effect on providing an unbiased  
353 representation of the controversial topic to the user. This is because these  
354 metrics were devised particularly for evaluating bias in the ranked outputs  
355 instead of SERPs.

356 Although the proposed measures by Yang & Stoyanovich (2017) are valuable  
357 in the context of measuring bias in ranked outputs where the individuals are  
358 being ranked and some of these individuals are the members of the protected  
359 group ( $g_1$ ), these measures have the aforementioned limitations. These limita-  
360 tions are particularly visible for content bias evaluation where the web docu-  
361 ments are being ranked by search engines in a typical IR setting. In this paper  
362 we address these limitations by proposing a family of fairness-aware measures  
363 with the main purpose of evaluating content bias in SERPs, based on standard  
364 utility-based IR evaluation measures.

365 Zehlike *et al.* (2017), based on Yang & Stoyanovich (2017)'s work, propose  
366 an algorithm to test the statistical significance of a fair ranking. Beutel *et al.*  
367 (2019) propose a pairwise fairness measure for recommender systems. However,  
368 the authors, unlike us, measure fairness on personalized recommendations and  
369 do not consider relevance, while we work in an unpersonalized information re-  
370 trieval setting and we do consider relevance. Kallus & Zhou (2019) investigate  
371 the fairness of predictive risk scores as a bipartite ranking task, where the

372 main goal is to rank positively labelled examples above negative ones. How-  
373 ever, their measures of bias based on the area under the ROC curve (AUC)  
374 are agnostic from the rank position at which a document has been retrieved.

### 375 2.3 Quantifying Search Engine Biases

376 Although the search engine algorithms are not transparent and available to  
377 external researchers, algorithm auditing techniques provide an effective means  
378 for systematically evaluating the results in a controlled environment (Sandvig  
379 *et al.*, 2014). Prior works leverage LDA-variant unsupervised methods and  
380 crowd-sourcing to analyse bias in content, or URL analysis for indexical bias.

381 Saez-Trumper *et al.* (2013) propose unsupervised methods to characterise  
382 different types of biases in online news media and in their social media commu-  
383 nities by also analysing political perspectives of the news sources. Yigit-Sert  
384 *et al.* (2016) investigate media bias by analysing the user comments along  
385 with the content of the online news articles to identify the latent aspects of  
386 two highly polarising topics in the Turkish political arena. Kulshrestha *et al.*  
387 (2017) quantify bias in social media by measuring the bias of the author of  
388 a tweet, while in Kulshrestha *et al.* (2018), bias in web search is quantified  
389 through a URL analysis for Google in political domain without any SERP  
390 content analysis. In our work, we consider the Google and Bing SERPs from  
391 news sources such as NY-Times, and BBC news in order to quantify bias  
392 through content analysis.

393 In addition to the unsupervised approaches, crowd-sourcing is a widely used  
394 mechanism to analyse bias in content. Crowd-sourcing is a common approach  
395 for labelling tasks in different research areas such as image & video annotation  
396 (Krishna *et al.*, 2017; Vondrick *et al.*, 2013), object detection (Su *et al.*, 2012),  
397 named entity recognition (Lawson *et al.*, 2010; Finin *et al.*, 2010), sentiment  
398 analysis (Räbiger *et al.*, 2018) and relevance evaluation (Alonso *et al.*, 2008;  
399 Alonso & Mizzaro, 2012). Yuen *et al.* (2011) provide a detailed survey of crowd-  
400 sourcing applications. As Yuen *et al.* (2011) suggest, crowd-sourcing can also  
401 be used for gathering opinions from the crowd. Mellebeek *et al.* (2010) use  
402 crowd-sourcing to classify Spanish consumer comments and show that non-  
403 expert Amazon Mechanical Turk (MTurk) annotations are viable and cost-  
404 effective alternative to expert ones. In this work, we use crowd-sourcing for  
405 collecting opinions of the public not about consumer products but controversial  
406 topics.

407 Apart from the content bias, there is another research area, namely index-  
408 ical bias. Indexical bias refers to the bias which is displayed in the selection of  
409 items, rather than in the content of retrieved documents, namely content bias  
410 (Mowshowitz & Kawaguchi, 2002b). Mowshowitz & Kawaguchi (2002a, 2005)  
411 quantify instead only indexical bias by using precision and recall measures.  
412 Moreover, the researchers approximate the *ideal* (i.e. norm) by the distribu-  
413 tion produced by a collection of search engines to measure bias. Yet, this  
414 may not be a *fair* bias evaluation procedure since the *ideal* itself should be

415 *unbiased*, whereas the SERPs of search engines may actually contain *bias*.  
416 Similarly, Chen & Yang (2006) use the same method in order to quantify  
417 indexical and content bias, however, content analysis was performed by repre-  
418 senting the SERPs with a weighted vector with different HTML tags without  
419 an in-depth analysis of the textual content. In this work, we evaluate content  
420 bias by analysing the textual contents of the Google and Bing SERPs, and  
421 we do not generate the *ideal* relying on the SERPs of other search engines in  
422 order to measure bias in a more *fair* way. In addition to the categorisation  
423 of the content and indexical bias analysis, prior methods used in auditing al-  
424 gorithms to quantify bias can also be divided into three main categories as  
425 *audience-based*, *content-based*, and *rater-based*. *Audience-based* measures fo-  
426 cus on identifying the political perspectives of media outlets and web pages  
427 by utilising the interests, ideologies, or political affiliations of its users, e.g.,  
428 likes and shares on Facebook (Bakshy *et al.*, 2015), based on the premise that  
429 readers follow the news sources that are closest to their ideological point of  
430 view (Mullainathan & Shleifer, 2005). Lahoti *et al.* (2018) model the problem  
431 of ideological leaning of social media users and media sources in the liberal-  
432 conservative ideology space on Twitter as a constrained non-negative matrix-  
433 factorisation problem. *Content-based* measures exploit linguistic features in  
434 textual content; Gentzkow & Shapiro (2010) extract frequent phrases of the  
435 different political partisans (Democrats, Republicans) from the Congress Re-  
436 ports. Then, the researchers come with the metric of media slant index to  
437 measure US newspapers’ political leaning. Finally, rater-based methods also  
438 exploit textual content and can be evaluated under the content-based methods.  
439 Unlike the content-based, the *rater-based* methods use ratings of people for the  
440 sentiment, partisan or ideological leaning of content instead of analysing the  
441 textual content linguistically. Rater-based methods generally leverage crowd-  
442 sourcing to collect the labels for the content analysis. For instance, Budak  
443 *et al.* (2016) quantify bias (partisanship) in US news outlets (newspapers and  
444 2 political blogs) for 15 selected queries related to a wide range of contro-  
445 versial issues about which Democrats and Republicans argue. The researchers  
446 use MTurk as a crowd-sourcing platform to obtain the topic and political slant  
447 labels, i.e. being positive towards Democrats or Republicans, of the articles.  
448 Similarly, Epstein & Robertson (2017) use crowd-sourcing to score individual  
449 search results and Diakopoulos *et al.* (2018) make use of the MTurk platform,  
450 i.e. rater-based approach, to get labels for the Google SERP websites by fo-  
451 cusing on the content and apply an audience-based approach through utilising  
452 the prior work of Bakshy *et al.* (2015) specifically for quantifying partisan  
453 bias. Our work follows a rater-based approach by making use of the MTurk  
454 platform for crowd-sourcing to analyse web search bias through stances and  
455 ideological leanings of the news articles instead of partisan bias in the textual  
456 contents of the SERPs.

457 There have been endeavors to audit partisan bias on web search. Diakopou-  
458 los *et al.* (2018) present four case studies on Google search results and to quan-  
459 tify partisan bias in the first page, they collect SERPs by issuing complete  
460 candidate names of the 2016 US presidential election as queries and utilise

crowd-sourcing to obtain the sentiment scores of the SERPs. They found that Google presented a higher proportion of negative articles for Republican candidates than the Democratic ones. Similarly, Epstein & Robertson (2017) present a case study for the election and use a browser extension to collect Google and Yahoo search data for the election-related queries, then use crowd-sourcing to score the SERPs. The researchers also found a left-leaning bias and Google was more biased than Yahoo. In their follow-up work, they found a small but significant ranking bias in the standard SERPs but not due to personalisation (Robertson *et al.*, 2018a). Similarly, researchers audit Google search after Donald Trump’s Presidential inauguration with a dynamic set of political queries using auto-complete suggestions (Robertson *et al.*, 2018b). Hu *et al.* (2019) conduct an algorithm audit and construct a specific lexicon of partisan cues for measuring political partisanship of Google Search snippets relative to the corresponding web pages. They define the corresponding difference as bias for this particular use case without making a robust search bias evaluation of SERPs from the user’s perspective. In this work, we introduce novel fairness-aware IR measures which involve rank information to evaluate content bias. For this, we use crowd-sourcing to obtain labels of the *news* SERPs returned towards the queries related to a wide-range of controversial topics instead of only political ones. With our robust bias evaluation measures, our main aim is to audit ideological bias in web search rather than solely partisan bias.

Apart from partisan bias, recent studies have investigated different types of bias for various purposes. Chen *et al.* (2018) investigate gender bias in the various resume search engines, which are platforms that help recruiters to search for suitable candidates and use statistical tests to examine two types of indirect discrimination: individual and group fairness. Similarly in another research study, authors investigate gender stereotypes by analyzing the gender distribution in image search results retrieved by Bing in four different regions (Otterbacher *et al.*, 2017). Researchers use the query of ‘person’ and the queries related to 68 character traits such as ‘intelligent person’, and the results show that photos of women are more often retrieved for ‘emotional’ and similar traits, whereas ‘rational’ and related traits are represented by photos of men. In a follow-up work, researchers conduct a controlled experiment via crowd-sourcing with participants from three different countries to detect bias in image search results (Otterbacher *et al.*, 2018). Demographic information along with measures of sexism are analysed together and the results confirm that sexist people are less likely to detect and report gender biases in the search results.

Raji & Buolamwini (2019) examine the impact of publicly naming biased performance results of commercial AI products in face recognition for directly challenging companies to change their products. Geyik *et al.* (2019) present a fairness-aware ranking framework to quantify bias with respect to protected attributes and improve the fairness for individuals without affecting the business metrics. The authors extended the metrics proposed by Yang & Stoyanovich (2017), of which we specified the limitations in Section 2.2, and evaluated their procedure using simulations with application to LinkedIn Tal-

ent Search. Vincent *et al.* (2019) measure the dependency of search engines on user-created content to respond to queries using Google search and Wikipedia articles. In another work, researchers propose a novel metric that involves users and their attention for auditing group fairness in ranked lists (Sapiezynski *et al.*, 2019). Gao & Shah (2019) propose a framework that effectively and efficiently estimate the solution space where fairness in IR is modelled as an optimisation problem with fairness constraint. Same researchers work on top-k diversity fairness ranking in terms of statistical parity and disparate impact fairness and propose entropy-based metrics to measure the topical diversity bias presented in SERPs of Google using clustering instead of a labelled dataset with group information (Gao & Shah, 2020). Unlike to their approach, our goal is to quantify search bias in SERPs rather than topical diversity. For this, we use a crowd-labelled dataset, thereby to evaluate bias from the user’s perspective with stance and ideological leanings of the documents.

In this context, we focus on proposing a new search bias evaluation procedure in ranked lists to quantify bias in the *news* SERPs. With the proposed robust fairness-aware IR measures, we also compare the relative bias of the two search engines through incorporating relevance and ranking information into the procedure without tracking the source of bias as discussed in Section 1. Our procedure can be used for the source of bias analysis as well which we leave as future work.

### 3 Search Engine Bias Evaluation Framework

In this section we describe our search bias evaluation framework. Then, we present the measures of bias and the proposed protocol to identify search bias.

#### 3.1 Preliminaries

Our first aim is to detect bias with respect to the distribution of stances expressed in the contents of the SERPs.

Let  $\mathcal{S}$  be the set of search engines and  $\mathcal{Q}$  be the set of queries about controversial topics. When a query  $q \in \mathcal{Q}$  is issued to a search engine  $s \in \mathcal{S}$ , the search engine  $s$  returns a SERP  $r$ . We define the stance of the  $i$ -th retrieved document  $r_i$  with respect to  $q$  as  $j(r_i)$ . A stance can have the following values: *pro*, *neutral*, *against*, *not-relevant*.

A document stance with respect to a topic can be:

- **pro** (👍) when the document is in favour of the controversial topic. The document describes more the pro aspects of the topic;
- **neutral** (👉) when the document does not support or help either side of the controversial topic. The document provides an impartial (fair) description about the pro and cons of the topic;
- **against** (👎) when the document is against the controversial topic. The document describes more the cons aspects of the topic;

547 – **not-relevant** (✘) when the document is not-relevant with respect to the  
 548 controversial topic.

549 For our analyses, we deliberately use recent *controversial* topics in US  
 550 that are the real debatable ones rather than the topics being possibly ex-  
 551 posed to false media balance, which occurs when the media present opposing  
 552 viewpoints as being more equal than the evidence supports, e.g. Flat Earth  
 553 debate (Grimes, 2016; Stokes, 2019). Our topic set contains abortion, illegal  
 554 immigration, gay marriage, and similar *controversial* topics which comprise  
 555 opposing points of view since complicated concepts concerning the identity,  
 556 religion, political or ideological leaning are the actual points where search en-  
 557 gines are more likely to provide biased results (Noble, 2018) and influence  
 558 people dramatically.

559 Our second aim is to detect bias with respect to the distribution of ideolog-  
 560 ical leanings expressed in the contents of the SERPs. We do this by associating  
 561 each query  $q \in \mathcal{Q}$  belonging to a controversial topic to one *current* ideological  
 562 leaning. Then, combining the stances for each  $r_i$  and the associated ideological  
 563 leaning of  $q$  we can measure the ideological bias of the content of a given SERP,  
 564 e.g. if a topic belongs to a specific ideology and a document retrieved for this  
 565 topic has a pro stance, we consider this document to be biased towards this  
 566 ideology. We define the ideological leaning of  $q$  as  $j(q)$ . An ideological leaning  
 567 can have the following values: *conservative, liberal, both or neither*.

568 A topic ideological leaning can be:








- 569 – **conservative** (●) when the topic is part of the conservative policies. The  
 570 conservatives are in favour of the topic;
- 571 – **liberal** (●) when the topic is part of the liberal policies. The liberals are  
 572 in favour of the topic;
- 573 – **both or neither** (○) when both or neither policies are either in favour or  
 574 against the topic.



575 For reference, Table 1 shows a summary of all the symbols, functions and  
 576 labels used in this paper.

### 577 3.2 Measures of Bias



578 Based on the aforementioned definition provided in Section 1, bias can be  
 579 quantified by measuring the degree of deviation of the distribution of doc-  
 580 uments from the *ideal* one. To give a broad definition of an ideal list poses  
 581 problems; but in the scope of this work for *controversial* topics, we can mention  
 582 the existence of bias in a ranked list retrieved by a search engine if the pre-  
 583 sented information *significantly deviates* from true likelihoods (White, 2013).  
 584 As justified in Section 1, in the scope of this work we focus on *equality of*  
 585 *output*, thus we accept the true likelihoods of different views as *equal* rather  
 586 than computing them from the corresponding base rates. Therefore using the  
 587 proposed definition reversely, we can assume that the *ideal* list is the one that  
 588 minimises the difference between two opposing views, which we indicate here

**Table 1** Symbols, functions, and labels used throughout the paper


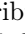
Symbols	
$\mathcal{S}$	set of search engines.
$s$	a search engine $s \in \mathcal{S}$ .
$\mathcal{Q}$	set of queries.
$q$	a query $q \in \mathcal{Q}$ .
$r$	a ranked list of the given SERP (list of retrieved documents).
$r_i$	the document in $r$ retrieved at rank $i$ .
$ r $	size of $r$ (number of documents in the ranked list).
$n$	number of documents considered in $r$ (cut-off).
Functions	
$j(r_i)$	returns the label associated to $r_i$ .
$f(r)$	an evaluation measure for SERPs.
Labels	
	pro stance.
	neutral stance.
	against stance.
	not-relevant stance.
	conservative ideological leaning.
	liberal ideological leaning.
	both or neither ideological leanings.

589 as  and  in the context of stances. Formally, we measure the *stance* bias  
 590 in a SERP  $r$  as follows:

$$591 \quad \beta_f(r) = f_{\text{thumbs up}}(r) - f_{\text{thumbs down}}(r), \quad (3)$$

592 where  $f$  is a function that measures the likelihood of  $r$  in satisfying the in-  
 593 formation need of the user about the view  and the view . We note that  
 594 *ideological* bias is measured in the same way by transforming the stances of  
 595 the documents into ideological leanings which will be explained in Section 4.2.  
 596 Before defining  $f$ , from Eq. (3), we define the mean bias (MB) of a search  
 597 engine  $s$  as:

$$598 \quad \text{MB}_f(s, \mathcal{Q}) = \frac{1}{|\mathcal{Q}|} \sum_{q \in \mathcal{Q}} \beta_f(s(q)).$$

599 An unbiased search engine would produce a mean bias of 0. A limitation of MB  
 600 is that if a search engine is biased towards the  view on one topic and bias  
 601 towards the  view on another topic, these two contributions will cancel each  
 602 other out. In order to avoid this limitation we also define the mean absolute  
 603 bias (MAB), which consists in taking the absolute value of the bias for each  
 604  $r$ . Formally, this is defined as follows:

$$605 \quad \text{MAB}_f(s, \mathcal{Q}) = \frac{1}{|\mathcal{Q}|} \sum_{q \in \mathcal{Q}} |\beta_f(s(q))|. \quad (4)$$

606 An unbiased search engine produces a mean absolute bias of 0. Although this  
 607 measure defined in Eq. (4) solves the limitation of MB, MAB says nothing

608 about towards which view the search engine is biased, making these two mea-  
 609 sures of bias complementary.

610 In IR the likelihood of  $r$  in satisfying the information need of users is  
 611 measured via retrieval evaluation measures. Among these measures we selected  
 612 3 *utility-based* evaluation measures. This class of evaluation measures quantify  
 613  $r$  in terms of its worth to the user and are normally computed as a sum of the  
 614 information gain summed over the relevant documents retrieved by  $r$ . The 3  
 615 IR evaluation measures used in the following experiments are: P@ $n$ , RBP, and  
 616 DCG@ $n$ .

617 P@ $n$  for the  $\heartsuit$  view is formalised as in Eq. (2). However, differently from  
 618 the previous definition of  $j(r_i)$  where the only possible outcomes are  $g_1$  and  $g_2$   
 619 for the document  $r_i$ , here  $j$  can return any of the label associated to a stance  
 620 ( $\heartsuit$ ,  $\spadesuit$ ,  $\clubsuit$ , and  $\times$ ). Hence, only pro and against documents, that are relevant to  
 621 the topic, are taken into account, since  $j(r_i)$  returns *neutral* and *not-relevant*  
 622 when otherwise. Substituting Eq. (2) to Eq. (3) we obtain the first measure of  
 623 bias:

$$624 \quad \beta_{P@n}(r) = \frac{1}{n} \sum_{i=1}^n ([j(r_i) = \heartsuit] - [j(r_i) = \spadesuit]).$$

625 The main limitation of this measure of bias is that it has a weak concept  
 626 of ranking, i.e. the first  $n$  documents contribute equally to the bias score.  
 627 The next two evaluation measures overcome this issue by defining discount  
 628 functions.

629 RBP weights every document based on the coefficients of a normalised  
 630 geometric series with value  $p \in ]0, 1[$ , where  $p$  is a parameter of RBP. Similarly  
 631 to what is done for P@ $n$ , we reformulate RBP to measure bias as follows:

$$632 \quad \text{RBP}_{\heartsuit} = (1 - p) \sum_{i=1}^n p^{i-1} [j(r_i) = \heartsuit]. \quad (5)$$

633 Substituting Eq. (5) to Eq. (3) we obtain:

$$634 \quad \beta_{\text{RBP}}(r) = (1 - p) \sum_{i=1}^n p^{i-1} ([j(r_i) = \heartsuit] - [j(r_i) = \spadesuit]).$$

635 DCG@ $n$ , instead, weights each document based on a logarithmic discount  
 636 function. Similarly to what is done for P@ $n$  and RBP, we reformulate DCG@ $n$   
 637 to measure bias as follows:

$$638 \quad \text{DCG}_{\heartsuit}@n = \sum_{i=1}^n \frac{1}{\log(i+1)} [j(r_i) = \heartsuit]. \quad (6)$$

639 Substituting Eq. (6) to Eq. (3) we obtain:

$$640 \quad \beta_{\text{DCG@}n}(r) = \sum_{i=1}^n \frac{1}{\log(i+1)} ([j(r_i) = \heartsuit] - [j(r_i) = \spadesuit]) \quad (7)$$

641 Since we are evaluating web-users, for P@ $n$  and DCG@ $n$  we set  $n = 10$   
 642 and for RBP we set  $p = 0.8$ . This last formulation (Eq. (7)), although it



643 looks similar to the rND measure, it does not suffer from the four limitations  
644 introduced in Section 2.2. In particular all these presented measures of bias: 1)  
645 do not focus on one group; 2) use a binary score associated to the document  
646 stance or ideological leaning, similar to the way these measures are used in IR  
647 when considering relevance; also like in IR 3) can be computed at each rank;  
648 4) exclude non-relevant documents from the measurement of bias and; the  
649 framework 5) provides various user models associated to the 3 IR evaluation  
650 measures:  $P@n$ ,  $DCG@n$ , and RBP.

### 651 3.3 Quantifying Bias

652 Using the measures of bias defined in the previous section we quantify the bias  
653 of the two search engines, Bing and Google using the *news versions* of these  
654 search engines. Then, we compare them thereof. Following, we describe each  
655 step of the proposed procedure used to quantify bias in SERPs.

- 656 – **News Articles in SERPs.** We obtained the controversial queries issued  
657 for searching from ProCong.org [2018] and applied some filtering steps on  
658 the initial query set. After filtering, the final query set size became 57. We  
659 submitted each query in the final query set to the US News search engines  
660 of Google and Bing using a US proxy. Then, we extracted the whole corpus  
661 returned by both engines in response to all the queries in the set. Note that  
662 the data collection process was done in a controlled environment such that  
663 the queries are sent to the search engines at the same time. For more details  
664 about the selection of the queries and crawling the SERPs, please refer to  
665 the previous phase of our analysis. After having crawled all the SERPs  
666 returned from both engines and extracted their contents, we annotated  
667 the top 10 documents. We obtained the stance label of each document  
668 with respect to the queries via crowd-sourcing. To label the ideological  
669 leaning of queries, we also used crowd-sourcing. To obtain the ideologies of  
670 documents, we transformed the stance labels into ideologies based on the  
671 ideological leaning of their corresponding queries. The details about our  
672 crowdsourcing campaigns as well as the transformation process can also be  
673 found in the first phase of our analysis.
- 674 – **Bias Evaluation.** We compute the bias measures for every SERP with  
675 all three IR-based measures of bias:  $P@n$ , RBP, and  $DCG@n$ . We then  
676 aggregate the results using the two measures of bias, MB and MAB.
- 677 – **Statistical Analysis.** To identify whether the bias measured is not a  
678 byproduct of randomness, we compute a one-sample t-test: the null hy-  
679 pothesis is that no difference exists and that the true mean is equal to  
680 zero. If this hypothesis is rejected, hence there is a significant difference  
681 and we claim that the evaluated search engine is biased. Then, we com-  
682 pare the difference in bias measured across the two search engines using a  
683 two-tailed paired t-test: the null hypothesis is that the difference between  
684 the two true means is equal to zero. If this hypothesis is rejected, hence

**Table 2** All controversial topics, topics marked with red dots are conservative and blue for liberal

• <b>Abortion:</b> Should Abortion Be Legal?	• <b>Alternative Energy vs. Fossil Fuels:</b> Can Alternative Energy Effectively Replace Fossil Fuels?	• <b>Animal Testing:</b> Should Animals Be Used for Scientific or Commercial Testing?
• <b>Banned Books:</b> Should Parents or Other Adults Be Able to Ban Books from Schools and Libraries?	• <b>Bill Clinton:</b> Was Bill Clinton a Good President?	• <b>Born Gay? Origins of Sexual Orientation:</b> Is Sexual Orientation Determined at Birth?
○ <b>Cell Phones Radiation:</b> Is Cell Phone Radiation Safe?	• <b>Climate Change:</b> Is Human Activity Primarily Responsible for Global Climate Change?	○ <b>College Education Worth It?:</b> Is a College Education Worth It?
• <b>Concealed Handguns:</b> Should Adults Have the Right to Carry a Concealed Handgun?	• <b>Corporal Punishment:</b> Should Corporal Punishment Be Used in K-12 Schools?	• <b>Corporate Tax Rate &amp; Jobs:</b> Does Lowering the Federal Corporate Income Tax Rate Create Jobs?
• <b>Cuba Embargo:</b> Should the United States Maintain Its Embargo against Cuba?	○ <b>Daylight Savings Time:</b> Should the United States Keep Daylight Saving Time?	○ <b>Drinking Age - Lower It?:</b> Should the Drinking Age Be Lowered from 21 to a Younger Age?
• <b>Drone Strikes Overseas:</b> Should the United States Continue Its Use of Drone Strikes Abroad?	○ <b>Drug Use in Sports:</b> Should Performance Enhancing Drugs (Such as Steroids) Be Accepted in Sports?	• <b>Electoral College:</b> Should the United States Use the Electoral College in Presidential Elections?
• <b>Euthanasia &amp; Assisted Suicide:</b> Should Euthanasia or Physician-Assisted Suicide Be Legal?	○ <b>Vaping E-Cigarettes:</b> Is Vaping with E-Cigarettes Safe?	• <b>Felon Voting:</b> Should Felons Who Have Completed Their Sentence (Incarceration, Probation, and Parole) Be Allowed to Vote?
○ <b>Fighting in Hockey:</b> Should Fighting Be Allowed in Hockey?	• <b>Gay Marriage:</b> Should Gay Marriage Be Legal?	○ <b>Gold Standard:</b> Should the United States Return to a Gold Standard?
○ <b>Golf - Is It a Sport?:</b> Is Golf a Sport?	• <b>Illegal Immigration:</b> Should the Government Allow Immigrants Who Are Here Illegally to Become US Citizens?	• <b>Israeli-Palestinian Two-State Solution:</b> Is a Two-State Solution (Israel and Palestine) an Acceptable Solution to the Israeli-Palestinian Conflict?
○ <b>Lowering the Voting Age to 16:</b> Should the Voting Age Be Lowered to 16?	• <b>Medical Marijuana:</b> Should Marijuana Be a Medical Option?	○ <b>Milk - Is It Healthy?:</b> Is Drinking Milk Healthy for Humans?
• <b>Minimum Wage:</b> Should the Federal Minimum Wage Be Increased?	• <b>National Anthem Protest:</b> Is Refusing to Stand for the National Anthem an Appropriate Form of Protest?	• <b>Net Neutrality:</b> Should Net Neutrality Be Restored?
• <b>Obamacare:</b> Obamacare Is the Patient Protection and Affordable Care Act (Obamacare) Good for America?	• <b>Obesity a Disease?:</b> Is Obesity a Disease?	○ <b>Olympics:</b> Are the Olympic Games an Overall Benefit for Their Host Countries and Cities?
○ <b>Penny - Keep It?:</b> Should the Penny Stay in Circulation?	○ <b>Police Body Cameras:</b> Should Police Officers Wear Body Cameras?	• <b>Prescription Drug Ads:</b> Should Prescription Drugs Be Advertised Directly to Consumers?
• <b>Prostitution - Legalize It?:</b> Should Prostitution Be Legal?	• <b>Right to Health Care:</b> Should All Americans Have the Right (Be Entitled) to Health Care?	• <b>Ronald Reagan:</b> Was Ronald Reagan a Good President?
• <b>Sanctuary Cities:</b> Should Sanctuary Cities Receive Federal Funding?	• <b>School Uniforms:</b> Should Students Have to Wear School Uniforms?	• <b>School Vouchers:</b> Are School Vouchers a Good Idea?
○ <b>Social Media:</b> Are Social Networking Sites Good for Our Society?	• <b>Social Security Privatization:</b> Should Social Security Be Privatized?	• <b>Standardized Tests:</b> Is the Use of Standardized Tests Improving Education in America?
• <b>Student Loan Debt:</b> Should Student Loan Debt Be Easier to Discharge in Bankruptcy?	○ <b>Tablets vs. Textbooks:</b> Should Tablets Replace Textbooks in K-12 Schools?	• <b>Teacher Tenure:</b> Should Teachers Get Tenure?
• <b>Under God in the Pledge:</b> Should the Words "Under God" Be in the US Pledge of Allegiance?	• <b>Universal Basic Income:</b> Is Universal Basic Income a Good Idea?	○ <b>Vaccines for Kids:</b> Should Any Vaccines Be Required for Children?
• <b>Vegetarianism:</b> Should People Become Vegetarian?	○ <b>Video Games and Violence:</b> Do Violent Video Games Contribute to Youth Violence?	○ <b>Voting Machines:</b> Do Electronic Voting Machines Improve the Voting Process?

685 there is a significant difference, we claim that there is a difference in bias  
686 between the two search engines.

687

## 688 4 Experimental Setup

689 In this section we provide a description of our experimental setup based on  
690 the proposed method as defined in Section 3.3.

## 691 4.1 Material

692 We obtained all the controversial topics from ProCon.org (2018). ProCon.org is  
693 a non-profit charitable organisation that provides an online resource for search  
694 on controversial topics. ProCon.org selects the topics that are controversial and  
695 important to many US citizens by also taking the readers' suggestions into  
696 account. We collected all 74 controversial topics with their topic questions  
697 from the website. Then, we applied three filters on these topics for practical  
698 reasons without deliberately selecting any topics. The first filter selects only  
699 the *polar* questions, also known as yes-no questions because they have no  
700 different sides for the analysis. This filter decreased the topic set size from  
701 74 to 70. The second filter removes the topics that do not contain up-to-date  
702 information in their topic pages provided by ProCon.org since they are not  
703 *recent* controversial topics and would not return up-to-date results. With the  
704 second filter, the number of topics became 64. Lastly, the third filter only  
705 includes the topics if both search engines return results for the corresponding  
706 topic questions, otherwise the comparison analysis would not be possible. After  
707 the last filter, the final topic set became the size of 57. Table 2 contains the  
708 full list of controversial topic titles with questions used in this study.

709 We used the topic questions of these 57 topics for crawling. For example,  
710 the topic question of the topic title 'abortion' is 'Should Abortion Be Legal?'.  
711 The topic questions reflect the main debate on the corresponding controversial  
712 topics and we used them as they are (i.e. including upper-cased characters,  
713 without removing punctuation, etc.) for querying the search engines.

714 We collected the news search results in *incognito mode* to avoid any per-  
715 sonalisation effect. Thus, the retrieved SERPs are not specific to anyone, but  
716 (presumably) general to US users. We submitted each topic question to US  
717 News search engines of Google and Bing using a US proxy. Since we used the  
718 *news versions* of the two search engines, sponsoring results which may affect  
719 our analysis did not appear in the news search results at all. Then, we firstly  
720 crawled the URLs of the retrieved results for the same topic question to min-  
721 imise the time lags between the search engines since the SERP of the same  
722 topic may vary over time. Subsequently, we extracted the textual contents  
723 of the top-10 documents using the crawled URLs. By this way, the time span  
724 between the SERPs of Google and Bing for each controversial topic (whole cor-  
725 pus) became 2-3 minutes on average. Moreover, before starting the crawling  
726 process, we firstly made some experiments with a small set of topics (different  
727 from the topic set provided in the paper) in the news search as well as default  
728 search and did not observe significant changes especially in the top-10 docu-  
729 ments of the news search even in 10-15 minutes time lags. This indicates that  
730 the news search is less dynamic than default search and we believe that the  
731 2-3 minutes of time lags would not drastically affect the search results.

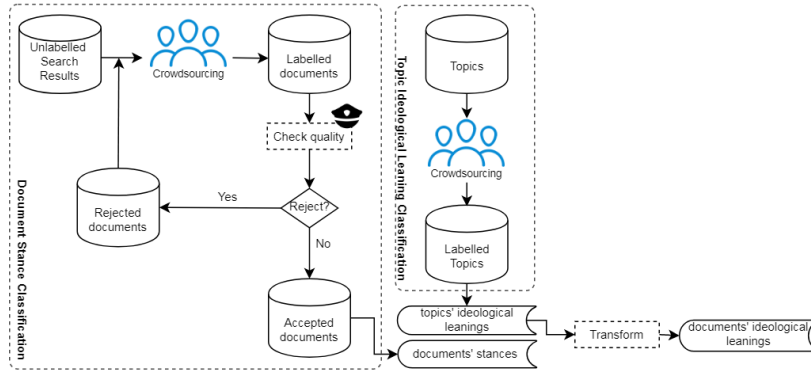


Fig. 1 Flow-chart of the crowd-sourcing campaigns

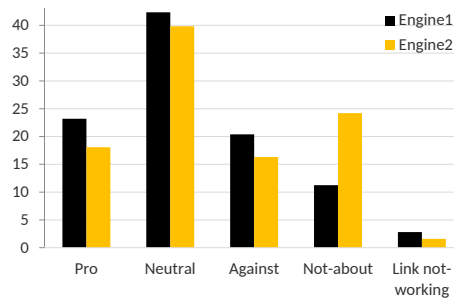
## 4.2 Crowd-sourcing Campaigns

The end-to-end process of obtaining stances and ideological leanings is shown in the flow-chart in Figure 1. The emphasised (dotted) parts of the flow-chart show the steps of the Document Stance Classification (DSC) and Topic Ideological Learning Classification (TILC).

The DSC process inputs unlabelled top-10 search results, crawled by the data collection procedure described in Section 4.1, and outputs the stance labels of all these documents via crowd-sourcing with respect to the topic questions ( $\mathcal{Q}$ ) used to retrieve them. As displayed in the flow-chart, the TILC process uses crowd-sourcing to output the ideological leanings of all topic questions ( $\mathcal{Q}$ ). Then, the accepted stance labels of all documents, acquired from the DSC process are transformed into ideological leaning labels based on the assigned ideology of their corresponding topic questions. The steps of obtaining document labels in stance and ideological leaning detection are described below.

To label the stance of each document with respect to the topic questions ( $\mathcal{Q}$ ) we used crowd-sourcing. We selected MTurk as a crowd-sourcing platform. In this platform, to obtain high quality crowd-labels task properties were set as follows. Since the topics are mostly related to US, we selected crowd-workers only from US. Moreover, we tried to find qualified and experienced workers by setting the following thresholds: Human Intelligence Task (HIT) approval rate percentage should be greater than 95% and number of HITs approved should be greater than 1000 for each worker. We set the wage as 0.15\$ and time allowed was 30 minutes per HIT. Each document was judged by three crowd-workers.

To classify the stance of a document we asked crowd-workers to label, given a controversial topic question, the stance of a document in pro, neutral, against, not-relevant, or link not-working. Before the task was assigned, instructions were given to a worker in three groups from general to specific. Initially, workers were provided an overview of the stance detection task, then



**Fig. 2** Percentages of the document stance labels annotated by crowd-workers

762 steps of the task were listed, i.e. read the topic question, open the news article  
 763 link etc., and finally, rules and tips were displayed. This last part contained  
 764 definitions of having a pro, neutral or against stance as given in Section 3.1  
 765 above. Additionally, we included a clue for workers saying that title of the ar-  
 766 ticle may give you a general idea about the stance, however it is not sufficient  
 767 to determine its overall viewpoint and then request workers to read also the  
 768 rest of the article. Apart from these, at the end of the page we put a warning  
 769 and informed the workers that some of the answers were known to us and  
 770 we may reject their HITs, i.e. single, self-contained task for a worker, based  
 771 on evaluation. Then, in the following page a HIT was shown to the worker  
 772 with a topic question (query), link to the news article whose stance will be  
 773 determined by repeating/reminding the main question of the stance detection  
 774 task.

775 In order to obtain reliable annotations, we first annotated a randomly cho-  
 776 sen set of documents later used to check the quality of crowd-labels as specified  
 777 in the warning to the workers. With these expert labels, we rejected low qual-  
 778 ity annotations and requested new labels for those documents. This iterative  
 779 process continued until we obtained all the document labels. At the end of this  
 780 iterative process, for the sake of label reliability, we computed two agreement  
 781 scores on the approved labels for document stance detection reported in Ta-  
 782 ble 3. The reported inter-rater agreement scores are the percent agreements  
 783 between the corresponding annotators. We looked at pairwise agreement; put  
 784 1 if there is an agreement and 0, otherwise. Then we computed the mean  
 785 for the fractions. Reported Kappa score for document stance classification is  
 786 considered *fair* agreement. Previously, researchers reported a Kappa score of  
 787 the inter-rater agreement between experts (0.385) instead of crowd-workers for  
 788 the same task, i.e. document stance classification in SERPs towards a different  
 789 query set which includes controversial topics as well as popular products, by  
 790 claiming MTurk workers had difficulty with the task (Alam & Downey, 2014).  
 791 Although our task seems to be more challenging, i.e. the queries are only about  
 792 controversial issues, our reported Kappa score for MTurk workers is compara-

**Table 3** Crowd-workers Agreement

Campaign	Inter-rater	Fleiss-Kappa
Document Stance	0.4968	0.3500
Topic Ideological Leaning	0.5281	0.3478

**Table 4** Performance of the search engines, p-values of a two-tailed paired t-test computed between engine 1 and 2

	P@10	RBP	DCG@10
Engine 1	0.8509	0.7708	3.9114
Engine 2	0.7404	0.6886	3.4773
p-value	< 0.001	< 0.001	< 0.01

793 ble to their expert agreement score, which we believe to be sufficient due to  
 794 the subjective nature and difficulty of the task.

795 The distribution of the accepted stance labels for the search results of each  
 796 search engine is displayed in Figure 2. One may argue that for a query about  
 797 a controversial topic issued to a news search engine, its SERP would mostly  
 798 contain controversial articles that support one dominant viewpoint towards a  
 799 given topic. Hence, informational pages or articles adequately discussing dif-  
 800 ferent viewpoints of the topic, i.e. documents that have a neutral stance, would  
 801 never get a chance to be included in the analysis. However, the distribution in  
 802 Figure 2 refutes this argument by showing that the majority of the labels for  
 803 both search engines is actually *neutral*.

804 To identify the ideological leaning of each topic, we again used crowd-  
 805 sourcing as displayed in Figure 1. We asked the crowd-workers to classify each  
 806 topic as: conservative, liberal, or both or neither. To get high quality annota-  
 807 tions also for topic ideology detection, worker properties were set as the same  
 808 with the stance detection. We again selected crowd-workers only from US. The  
 809 wage per HIT was set as 0.1\$ and the time allowed was 5 minutes. Similarly  
 810 to the stance detection, in the informational page we gave an overview, listed  
 811 the steps and lastly provided the rules & tips. For this task, last part con-  
 812 tained the ideological leaning definitions as given in Section 3.1. Additionally,  
 813 we requested the workers to evaluate the ideological leaning of a given topic  
 814 based on the current ideological climate and warned them related to the re-  
 815 jection of their HITs as before. In the next page, the workers were shown a  
 816 HIT with a topic question (query), i.e. one of the main debates of the corre-  
 817 sponding topic, and asked the worker the following: *Which ideological group*  
 818 *would answer favourably to this question?*. The topics assigned to conservative  
 819 or liberal leanings have been decided based on the judgment of five annotators  
 820 with majority-voting. The leanings of the topics are shown in Table 2. Two  
 821 agreement scores computed on the judgments for ideological leaning detection  
 822 are also reported in Table 3.

823 To map the stance from the *pro-to-against* to the *conservative-to-liberal*,  
 824 we applied a simple transformation to the documents. This transformation is

**Table 5** Stance bias of the search engines, p-values of a two-tailed paired t-test computed between engine 1 and 2

		P@10	RBP	DCG@10
MB	Engine 1	0.0281	0.0197	0.1069
	Engine 2	0.0175	0.0271	0.1142
	p-value	> 0.05	> 0.05	> 0.05
MAB	Engine 1	0.2596	0.2738	1.3380
	Engine 2	0.2246	0.2266	1.0789
	p-value	> 0.05	> 0.05	> 0.05

**Table 6** Ideological bias of the search engines, p-values of a two-tailed paired t-test computed between engine 1 and 2

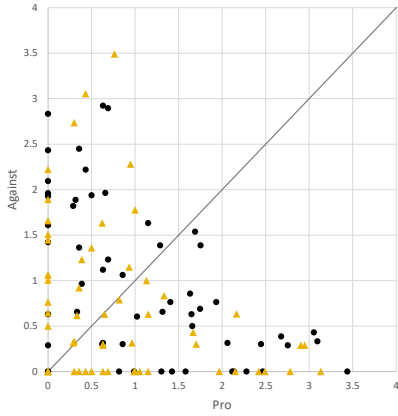
		P@10	RBP	DCG@10
MB	Engine 1	-0.1368	-0.1247	-0.6290
	Engine 2	-0.1289	-0.1386	-0.6591
	p-value	> 0.05	> 0.05	> 0.05
MAB	Engine 1	0.2579	0.2894	1.3989
	Engine 2	0.2184	0.2158	1.0456
	p-value	> 0.05	< 0.05	< 0.05

825 needed because there may be documents which have a pro stance, for example,  
826 towards *abortion* and *Cuba embargo*. Though these documents have the same  
827 stance, they have different ideological leanings since having a pro stance on  
828 *abortion* implies a *liberal leaning*, whereas a pro stance on *Cuba embargo* im-  
829 plies a *conservative leaning*. For some topics (as in the case of *Cuba embargo*),  
830 we can directly interpret the *pro-to-against* stance labels of search results as  
831 *conservative-to-liberal* ideological leaning labels while for other topics (as in  
832 the case of the *abortion*) as *liberal-to-conservative*. On the other hand, for  
833 those topics such as *vaccines for kids*, which crowded label resulted in both or  
834 neither, the conservative-to-liberal or liberal-to-conservative transformation  
835 was not meaningful and therefore eliminated by our analysis. We note that  
836 within budget constraints, the crowd-sourcing protocol was designed to ob-  
837 tain crowd-labels with high-quality by labelling (expert) the random sample  
838 of documents, applying iterative process and majority voting on these labels.

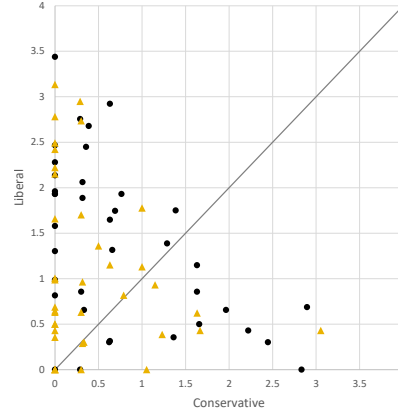
### 839 4.3 Results

840 In Table 4 we present the performance of the two search engines. This is mea-  
841 sured over all the topics. A document is considered relevant when classified as  
842 pro, against, or neutral. The difference for all evaluation measures is statisti-  
843 cally significant.

844 In Table 5 we present the stance bias of the search engines. Note that  
845 for all the three measures of bias, P@10, RBP and DCG@10, lower value  
846 is better which means lower bias in the scope of this work as opposed to  
847 their corresponding classic IR measures. All MB and MAB scores are positive



**Fig. 3**  $DCG_{\blacktriangle}@10$  against  $DCG_{\blacklozenge}@10$  measured on stances – black points for engine 1 and yellow points for engine 2



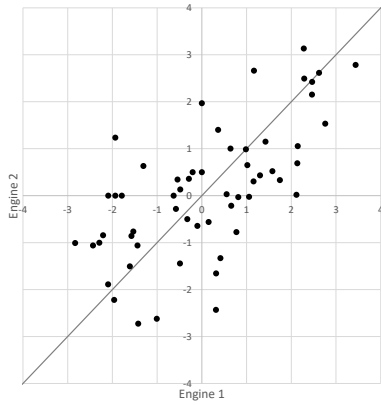
**Fig. 4**  $DCG_{\bullet}@10$  against  $DCG_{\bullet}@10$  measured on ideological leanings – black points for engine 1 and yellow points for engine 2

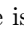

848 for all three IR evaluation measures. Also, the differences between the two  
 849 search engines for both MB and MAB measures are statistically not significant  
 850 and it is shown with the two-tailed pair t-test on these measures. In Table  
 851 6 we show the ideological bias. Similarly to Table 5, lower is better since  
 852 we use the same measures of bias. This table is similar to Table 5. Unlike  
 853 the Table 5, all MB scores are negative while all MAB scores are positive  
 854 for all three IR evaluation measures. The two-tailed paired t-test computed  
 855 on MBs to compare the difference in bias between engine 1 and engine 2,  
 856 this is statistically not significant. Nonetheless, the two-tailed test on MABs  
 857 is statistically not significant for the measure P@10; but it is statistically  
 858 significant for the measures RBP and DCG@10.

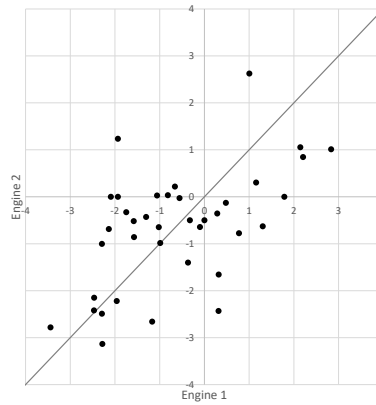
859 In Figure 3 we show how the topic-wise SERPs distribute over the pro-  
 860 against stance space for the measure DCG@10. The x-axis is the pro stance  
 861 score ( $DCG_{\blacktriangle}@10$ ) and the y-axis is the against stance score ( $DCG_{\blacklozenge}@10$ ).  
 862 Each point corresponds to the overall SERP score of a topic. Black points are  
 863 those SERPs retrieved by engine 1 and yellow points are those retrieved by  
 864 engine 2.

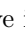

865 In Figure 5 we compare the overall stance bias score ( $\beta_{DCG@10}$ ), i.e. dif-  
 866 ference between the pro and against stance scores, of SERPs for each topic  
 867 measured on the two search engines. The x-axis is engine 1 and the y-axis is  
 868 engine 2. The points in positive coordinates denote the topics whose SERPs  
 869 are overall biased towards the pro stance, negative coordinates are for the  
 870 against stance.





**Fig. 5**  $\beta_{DCG@10}$  measured on stances, where positive is  and negative is 



**Fig. 6**  $\beta_{DCG@10}$  measured on leanings, where positive is  and negative is 

871 Figure 4 and Figure 6 are similar to Figure 3 and Figure 5 but instead of  
 872 measuring the stance bias we measure the ideological bias in the former case.  
 873 Therefore, Figure 4 displays how the overall SERPs of topics distribute over the  
 874 conservative-liberal ideological space for the measure DCG@10. Similarly, in  
 875 Figure 6 we compare the overall ideological bias score ( $\beta_{DCG@10}$ ), i.e. difference  
 876 between the conservative and liberal leaning scores, of the SERPs where the  
 877 points in positive coordinates stand for the topics that are biased towards the  
 878 conservative leaning, negative coordinates are for the liberal.

## 879 5 Discussion

880 Before investigating the existence of bias in SERPs, we initially compared the  
 881 retrieval performances of two search engines. In Table 4 we observe that the  
 882 performance of the two search engines is high but engine 1 is better than  
 883 engine 2 – their difference is statistically significant. This is verified across all  
 884 three IR evaluation measures.

885 Next, we verify if the search engines return biased results in terms of docu-  
 886 ment stances (RQ1) and if so, we further investigate if the engines suffer from  
 887 the same level of bias (RQ2) that the difference between the engines are not  
 888 statistically significant. In Table 5 all MB scores are positive and regarding the  
 889 RQ1, the engines seem to be biased towards the pro stance. We applied the  
 890 one-sample t-test on MB scores to check the existence of stance bias, i.e. if the  
 891 true mean is different from zero, as mentioned in Section 3.3. However, these  
 892 biases are statistically not significant which means that this expectation may  
 893 be the result of noise – there is not a systematic stance bias, i.e. preference of  
 894 one stance with respect to the other. Based on MAB scores, we can observe  
 895 that both engines suffer from an absolute bias. However, the difference between

the two engines is shown to be non-significant with the two-tailed t-test. These results show that both search engines are not biased towards a *specific* stance in returning results since there is no statistically significant difference from the *ideal* distribution. Nonetheless, for both engines there exists an absolute bias which can be interpreted as the expected bias for a topic question. These empirical findings imply that the search engines are biased for some topics towards the pro stance and for others towards the against stance.

The results are displayed in Figure 3. This figure refers to the values used to compute the MAB score of the DCG@10 column. It shows that the difference between the pro and against stances of both engines for topics is uniformly distributed. To note that, no topic can be located on the up-right area of the plot because the sum of their coordinates is bounded by the maximum possible DCG@10 score. Moreover we observe that topics are distributed similarly across the engines. This is also confirmed by Figure 5 where we can observe that the stance bias scores ( $\beta_{DCG@10}$ ), i.e. the differences between DCG@10 scores for the pro stance and DCG@10 scores for the against stance, of topics are somehow balanced between the up-right quadrant and the low-left quadrant. Moreover, these two quadrants are the area of agreement in stance between the two engines. The other two quadrants contain those topics where the engines disagree. Here we can conclude that the engines agree with each other in the majority of cases.

Lastly, we investigate if the search engines are biased in the ideology space (RQ3). Looking at MB scores in Table 6 we observe that both search engines seem to be biased towards the same ideological leaning – liberal (all MB scores are negative). Unlike the stance bias, one sample t-test on MB scores show that these expectations are statistically significant with different confidence values, i.e. p-value < 0.005 across all three IR measures for engine 2; whereas the same confidence value on P@10 for engine 1 and p-value < 0.05 on RBP and DCG@10. These results indicate that both search engines are biased towards the same leaning which is liberal. Comparing the two search engines on MB scores, we observe that their differences are statistically not significant, which means that the observed difference may be the result of random noise. Based on MAB, since all MAB scores are positive we can also observe that both engines suffer from an absolute bias. However, in contrast with what observed for the stance bias, this time there is a difference in expected ideological bias between the two search engines. For RBP and DCG@10 the difference between the engines is statistically significant. This finding and the different user models that these evaluation measures model suggest that the perceived bias by the users may change based on their behaviour. A user that always inspects the first 10 results (as modelled by P@10) may perceive the same ideological bias between engine 1 and engine 2, while a less systematic user, which just inspects the top results, may perceive that engine 1 is more biased than engine 2. Moreover, comparing this finding with the performance of the engines, we can observe that the better performing engine is more biased than the worse performing one.

941 Comparing Figure 4 with Figure 3 we observe that in Figure 4 the points  
942 look less uniformly distributed than in Figure 3. Topics are mostly on the  
943 liberal side. Moreover, engine 2 has fewer points on the conservative side than  
944 engine 1. Comparing Figure 6 with Figure 5, we observe that the engines in  
945 Figure 6 are more biased towards the liberal side with respect to what observed  
946 in Figure 5. Also, we observe that the engines mostly agree – most of the points  
947 are placed on the up-right and low-left quadrants.

948 In conclusion, we find important to point out that it is not in the scope  
949 of this work to find the source of bias. As discussed in the introduction,  
950 bias may be a result of the input data, which may contain biases, or the  
951 search algorithm, which contains sophisticated features and specifically cho-  
952 sen algorithms that, although designed to be effective in satisfying information  
953 needs, may produce systematic biases. Nonetheless, we look at the problem  
954 from the user perspective and no matter where the bias comes from; the re-  
955 sults are biased as described. Our findings seem to be consistent with prior  
956 works (Epstein & Robertson, 2017; Diakopoulos *et al.*, 2018) that there exists  
957 liberal (left-leaning) partisan bias in SERPs; even in unpersonalised search  
958 settings (Robertson *et al.*, 2018a).

## 959 6 Limitations

960 This work has potential limitations. As stated in the introduction, we focus  
961 on a particular kind of bias, known as *statistical parity*, or more generally  
962 known as *equality of outcome* instead of *equality of opportunity* which uses  
963 query-specific base rates. In the context of the controversial topics where the  
964 document labels were obtained via crowd-sourcing, this bias measure, i.e. re-  
965 quiring equal representation of stances instead of query-specific base rates,  
966 made our experiments feasible. This is firstly because, not all of the query  
967 questions in our list have certain answers based on scientific facts, i.e. some  
968 of them are subjective queries. In investigating the equality of opportunity,  
969 queries can be further categorized as subjective and objective on top of our  
970 evaluation framework. For the objective queries, expert labels can be obtained  
971 and used as base rates, then search results can be evaluated by taking into  
972 account these base rates. Please note that our evaluation framework could  
973 better be applied to the controversial queries from the public’s perspective  
974 mainly where the goal is to have balanced SERPs instead of skewed results.  
975 We believe that some queries should be handled with a different framework  
976 since those queries are not intrinsically controversial such as *Is Holocaust real?*  
977 - there is only one correct answer without the need of a discussion.

978 Besides, the identification of the stance for the full ranking list is currently  
979 too expensive to get annotated via crowd-sourcing. To tackle this issue, a  
980 machine learning model can help us to automate the process of obtaining the  
981 stance labels. Another potential limitation is that some queries may not be  
982 real user queries. Nonetheless, we extracted the queries directly from their  
983 topic pages of the ProCon.org (2018) along with the topics. We deliberately

984 did not change the queries to avoid any interference/bias from our side on  
985 the results. In this work, we did not make a domain-specific selection of the  
986 topics, or apply any filtering as subjective/objective, rather we accepted them  
987 as *controversial* topics from the general public’s perspective which is the main  
988 scope of this work.

989 Apart from these, crowd-workers’ own personal biases may affect the la-  
990 belling process. For this reason, we tried to mitigate these biases by i. asking  
991 the workers to annotate stances rather than ideologies to make their judgment  
992 more objective, and ii. aggregating the final judgment coming from multiple  
993 workers. Additionally, our analysis refers to a specific point in time where the  
994 data was collected. To enable reproducibility and an easier comparison of these  
995 results at some point in the future, we made our dataset publicly available.  
996 Lastly, we note that this bias analysis can only be used as an indicator of po-  
997 tentially biased ranking algorithms because it is not enough in order to track  
998 the source of bias. In the scope of this work, we did not investigate the source  
999 of bias that may come from the data (input bias) or from the ranking mech-  
1000 anism (algorithmic bias) of the corresponding search engines. Despite these  
1001 potential limitations, we believe that our work is a good attempt to evaluate  
1002 bias in search results with new bias measures and a dataset crawled specifi-  
1003 cally for the search bias evaluation. Since the bias analysis is very complex, we  
1004 deliberately limited our scope and only focused on the bias analysis of *recent*  
1005 controversial topics in *news* search. Nonetheless, all these limitations lead us  
1006 to numerous interesting future directions.

## 1007 7 Conclusion & Future Work

1008 In this work we introduced new bias evaluation measures and a generalisable  
1009 evaluation framework to address the issue of web search bias in news search  
1010 results. We applied the proposed framework to measure stance and ideological  
1011 bias in the SERPs of Bing and Google as well as compare their relative bias  
1012 towards controversial topics. Our initial results show that both search engines  
1013 seem to be unbiased when considering the document stances and *ideologically*  
1014 biased when considering the document ideological leanings. In this work, we  
1015 intended to analyse SERPs without the effect of personalisation. Thus, these  
1016 results highlight that search biases exist even though the personalization ef-  
1017 fect is minimized and that search engines can empower users by being more  
1018 accountable.

1019 In the scope of this work we did not investigate the source of bias which we  
1020 left as future work, therefore the results can be seen as a potential indicator.  
1021 In our experiments, we gathered document stances via crowd-sourcing. Thus,  
1022 the obvious future work in this direction is to use automatic stance detection  
1023 methods instead of crowd-sourcing to obtain the document labels, thereby  
1024 evaluating bias in the whole corpus of retrieved SERPs to track the source of  
1025 bias. Moreover, investigating the workers’ bias in a follow-up work would be  
1026 interesting since it is very difficult to remove all biases in practice. In this work,

1027 we focus on *equality of outcome*; but using another bias measure, *equality of*  
1028 *opportunity* which takes into account the corresponding group proportions, i.e.  
1029 query-specific base rates, in the population would be an alternative follow-up  
1030 work. We plan to categorize queries as subjective and objective, then modify  
1031 the *ideal* ranking definition specifically for the objective queries based on the  
1032 corpus distributions. The bias analysis for the objective queries, particularly  
1033 the ones related to the critical domains such as health search, can be investi-  
1034 gated further on top of our evaluation framework which we believe to be an  
1035 interesting follow-up work. Furthermore, we plan to study the effect of local-  
1036 ization and personalization, i.e. how much the stances and ideological leanings  
1037 varied across users or the echo chamber effect, on SERPs, then incorporate  
1038 that study into our bias evaluation framework in the future.

### 1039 Compliance with Ethical Standards

1040 Author Emine Yilmaz previously worked as a research consultant for Microsoft  
1041 Research and she is currently a research consultant for Amazon Research.

### 1042 Acknowledgements

1043 We thank the reviewers for their comments. This work has been funded by the  
1044 EPSRC Fellowship titled "Task Based Information Retrieval", grant reference  
1045 number EP/P024289/1 and the visiting researcher programme of The Alan  
1046 Turing Institute.

### 1047 References

- 1048 (2018). Internetlivestats. <http://www.internetlivestats.com/>, accessed:  
1049 2018-10-06.
- 1050 (2018). Procon.org, procon.org - pros and cons of controversial issues. [https:](https://www.procon.org/)  
1051 [//www.procon.org/](https://www.procon.org/), accessed: 2018-07-31.
- 1052 (2018). Search engine statistics 2018. [https://www.smartinsights.com/](https://www.smartinsights.com/search-engine-marketing/search-engine-statistics/)  
1053 [search-engine-marketing/search-engine-statistics/](https://www.smartinsights.com/search-engine-marketing/search-engine-statistics/), accessed:  
1054 2018-10-06.
- 1055 99Firms (2019). Search engine statistics. [https://99firms.com/blog/](https://99firms.com/blog/search-engine-statistics/#gref)  
1056 [search-engine-statistics/#gref](https://99firms.com/blog/search-engine-statistics/#gref), accessed: 2019-09-06.
- 1057 Aktolga, E. & Allan, J. (2013). Sentiment diversification with different biases.  
1058 *Proceedings of the 36th international ACM SIGIR conference on Research*  
1059 *and development in information retrieval*, pp. 593–602, ACM.
- 1060 Alam, M.A. & Downey, D. (2014). Analyzing the content emphasis of web  
1061 search engines. *Proceedings of the 37th international ACM SIGIR conference*  
1062 *on Research & development in information retrieval*, pp. 1083–1086, ACM.
- 1063 Alonso, O. & Mizzaro, S. (2012). Using crowdsourcing for trec relevance as-  
1064 sessment. *Information processing & management*, **48**, 1053–1066.

- 1065 Alonso, O., Rose, D.E. & Stewart, B. (2008). Crowdsourcing for relevance  
1066 evaluation. *SIGIR forum*, vol. 42, pp. 9–15.
- 1067 Baeza-Yates, R. (2016). Data and algorithmic bias in the web. *Proceedings of*  
1068 *the 8th ACM Conference on Web Science*, pp. 1–1, ACM.
- 1069 Bakshy, E., Messing, S. & Adamic, L.A. (2015). Exposure to ideologically  
1070 diverse news and opinion on facebook. *Science*, **348**, 1130–1132.
- 1071 Bargh, J.A., Gollwitzer, P.M., Lee-Chai, A., Barndollar, K. & Trötschel, R.  
1072 (2001). The automated will: nonconscious activation and pursuit of behav-  
1073 ioral goals. *Journal of personality and social psychology*, **81**, 1014.
- 1074 Beutel, A., Chen, J., Doshi, T., Qian, H., Wei, L., Wu, Y., Heldt, L., Zhao,  
1075 Z., Hong, L., Chi, E.H. *et al.* (2019). Fairness in recommendation ranking  
1076 through pairwise comparisons. *arXiv preprint arXiv:1903.00780*, .
- 1077 Budak, C., Goel, S. & Rao, J.M. (2016). Fair and balanced? quantifying media  
1078 bias through crowdsourced content analysis. *Public Opinion Quarterly*, **80**,  
1079 250–271.
- 1080 Chelaru, S., Altingovde, I.S. & Siersdorfer, S. (2012). Analyzing the polarity  
1081 of opinionated queries. *European Conference on Information Retrieval*, pp.  
1082 463–467, Springer.
- 1083 Chelaru, S., Altingovde, I.S., Siersdorfer, S. & Nejd, W. (2013). Analyzing,  
1084 detecting, and exploiting sentiment in web queries. *ACM Transactions on*  
1085 *the Web (TWEB)*, **8**, 6.
- 1086 Chen, L., Ma, R., Hannák, A. & Wilson, C. (2018). Investigating the impact  
1087 of gender on rank in resume search engines. *Proceedings of the 2018 chi*  
1088 *conference on human factors in computing systems*, pp. 1–14.
- 1089 Chen, X. & Yang, C.Z. (2006). Position paper: A study of web search engine  
1090 bias and its assessment. *IW3C2 WWW*, .
- 1091 Culpepper, J.S., Diaz, F. & Smucker, M.D. (2018). Research frontiers in infor-  
1092 mation retrieval: Report from the third strategic workshop on information  
1093 retrieval in lorne (swirl 2018). *ACM SIGIR Forum*, vol. 52, pp. 46–47, ACM  
1094 New York, NY, USA.
- 1095 Demartini, G. & Siersdorfer, S. (2010). Dear search engine: what’s your opinion  
1096 about...?: sentiment analysis for semantic enrichment of web search results.  
1097 *Proceedings of the 3rd International Semantic Search Workshop*, p. 4, ACM.
- 1098 Diakopoulos, N., Trielli, D., Stark, J. & Mussenden, S. (2018). I vote for—how  
1099 search informs our choice of candidate. *Digital Dominance: The Power of*  
1100 *Google, Amazon, Facebook, and Apple*, M. Moore and D. Tambini (Eds.),  
1101 **22**.
- 1102 Diaz, A. (2008). Through the google goggles: Sociopolitical bias in search en-  
1103 gine design. *Web search*, pp. 11–34, Springer.
- 1104 Dutton, W.H., Blank, G. & Groselj, D. (2013). *Cultures of the internet: the*  
1105 *internet in Britain: Oxford Internet Survey 2013 Report*. Oxford Internet  
1106 Institute.
- 1107 Dutton, W.H., Reisdorf, B., Dubois, E. & Blank, G. (2017). Search and politics:  
1108 The uses and impacts of search in britain, france, germany, italy, poland,  
1109 spain, and the united states, .

- 1110 Elisa Shearer, K.E.M. (2018). News use across social media  
1111 platforms 2018. [https://www.journalism.org/2018/09/10/  
1112 news-use-across-social-media-platforms-2018/](https://www.journalism.org/2018/09/10/news-use-across-social-media-platforms-2018/).
- 1113 Epstein, R. & Robertson, R.E. (2015). The search engine manipulation effect  
1114 (seme) and its possible impact on the outcomes of elections. *Proceedings of  
1115 the National Academy of Sciences*, **112**, E4512–E4521.
- 1116 Epstein, R. & Robertson, R.E. (2017). A method for detecting bias in search  
1117 rankings, with evidence of systematic bias related to the 2016 presidential  
1118 election. *Technical Report White Paper no. WP-17-02*, .
- 1119 Epstein, R., Robertson, R.E., Lazer, D. & Wilson, C. (2017). Suppressing the  
1120 search engine manipulation effect (seme). *Proceedings of the ACM: Human-  
1121 Computer Interaction*, **1**, 42.
- 1122 Fang, Y., Si, L., Somasundaram, N. & Yu, Z. (2012). Mining contrastive opin-  
1123 ions on political texts using cross-perspective topic model. *Proceedings of  
1124 the fifth ACM international conference on Web search and data mining*, pp.  
1125 63–72, ACM.
- 1126 Finin, T., Murnane, W., Karandikar, A., Keller, N., Martineau, J. & Dredze,  
1127 M. (2010). Annotating named entities in twitter data with crowdsourcing.  
1128 *Proceedings of the NAACL HLT 2010 Workshop on Creating Speech and  
1129 Language Data with Amazon’s Mechanical Turk*, pp. 80–88, Association for  
1130 Computational Linguistics.
- 1131 Gao, R. & Shah, C. (2019). How fair can we go: Detecting the boundaries of  
1132 fairness optimization in information retrieval. *Proceedings of the 2019 ACM  
1133 SIGIR International Conference on Theory of Information Retrieval*, pp.  
1134 229–236.
- 1135 Gao, R. & Shah, C. (2020). Toward creating a fairer ranking in search engine  
1136 results. *Information Processing & Management*, **57**, 102138.
- 1137 Gentzkow, M. & Shapiro, J.M. (2010). What drives media slant? evidence from  
1138 us daily newspapers. *Econometrica*, **78**, 35–71.
- 1139 Geyik, S.C., Ambler, S. & Kenthapadi, K. (2019). Fairness-aware ranking in  
1140 search & recommendation systems with application to linkedin talent search.  
1141 *Proceedings of the 25th ACM SIGKDD International Conference on Knowl-  
1142 edge Discovery & Data Mining*, pp. 2221–2231.
- 1143 Ginger, G. & David, S. (2018). Google responds to trump,  
1144 says no political motive in search results. [https://  
1145 www.reuters.com/article/us-usa-trump-tech-alphabet/  
1146 google-responds-to-trump-says-no-political-motive-in-search-results-idUSKCN1LD1QP](https://www.reuters.com/article/us-usa-trump-tech-alphabet/google-responds-to-trump-says-no-political-motive-in-search-results-idUSKCN1LD1QP),  
1147 accessed: 2018-10-06.
- 1148 Goldman, E. (2008). Search engine bias and the demise of search engine utopi-  
1149 anism. *Web Search*, pp. 121–133, Springer.
- 1150 Grimes, D.R. (2016). Impartial journalism is laudable. but false balance is  
1151 dangerous. [https://www.theguardian.com/science/blog/2016/nov/08/  
1152 impartial-journalism-is-laudable-but-false-balance-is-dangerous](https://www.theguardian.com/science/blog/2016/nov/08/impartial-journalism-is-laudable-but-false-balance-is-dangerous),  
1153 accessed: 2019-08-15.
- 1154 Hu, D., Jiang, S., E. Robertson, R. & Wilson, C. (2019). Auditing the par-  
1155 tisanship of google search snippets. *The World Wide Web Conference*, pp.

693–704.

- 1156  
1157 Institute, A.P. (2014). The personal news cycle: How americans choose to get  
1158 their news. *American Press Institute*, .
- 1159 Kallus, N. & Zhou, A. (2019). The fairness of risk scores beyond classification:  
1160 Bipartite ranking and the xauc metric. *arXiv preprint arXiv:1902.05826*, .
- 1161 Krishna, R., Zhu, Y., Groth, O., Johnson, J., Hata, K., Kravitz, J., Chen, S.,  
1162 Kalantidis, Y., Li, L.J., Shamma, D.A. *et al.* (2017). Visual genome: Con-  
1163 necting language and vision using crowdsourced dense image annotations.  
1164 *International Journal of Computer Vision*, **123**, 32–73.
- 1165 Kulshrestha, J., Eslami, M., Messias, J., Zafar, M.B., Ghosh, S., Gummadi,  
1166 K.P. & Karahalios, K. (2017). Quantifying search bias: Investigating sources  
1167 of bias for political searches in social media. *Proceedings of the 2017 ACM*  
1168 *Conference on Computer Supported Cooperative Work and Social Comput-*  
1169 *ing*, pp. 417–432, ACM.
- 1170 Kulshrestha, J., Eslami, M., Messias, J., Zafar, M.B., Ghosh, S., Gummadi,  
1171 K.P. & Karahalios, K. (2018). Search bias quantification: investigating po-  
1172 litical bias in social media and web search. *Information Retrieval Journal*,  
1173 pp. 1–40.
- 1174 Lahoti, P., Garimella, K. & Gionis, A. (2018). Joint non-negative matrix  
1175 factorization for learning ideological leaning on twitter. *Proceedings of the*  
1176 *Eleventh ACM International Conference on Web Search and Data Mining*,  
1177 pp. 351–359.
- 1178 Lawson, N., Eustice, K., Perkowski, M. & Yetisgen-Yildiz, M. (2010). Anno-  
1179 tating large email datasets for named entity recognition with mechanical  
1180 turk. *Proceedings of the NAACL HLT 2010 workshop on creating speech*  
1181 *and language data with Amazon’s Mechanical Turk*, pp. 71–79, Association  
1182 for Computational Linguistics.
- 1183 Mellebeek, B., Benavent, F., Grivolla, J., Codina, J., Costa-Jussa, M.R. &  
1184 Banchs, R. (2010). Opinion mining of spanish customer comments with non-  
1185 expert annotations on mechanical turk. *Proceedings of the NAACL HLT*  
1186 *2010 workshop on Creating speech and language data with Amazon’s me-*  
1187 *chanical turk*, pp. 114–121, Association for Computational Linguistics.
- 1188 Mowshowitz, A. & Kawaguchi, A. (2002a). Assessing bias in search engines.  
1189 *Information Processing & Management*, **38**, 141–156.
- 1190 Mowshowitz, A. & Kawaguchi, A. (2002b). Bias on the web. *Communications*  
1191 *of the ACM*, **45**, 56–60.
- 1192 Mowshowitz, A. & Kawaguchi, A. (2005). Measuring search engine bias. *In-*  
1193 *formation processing & management*, **41**, 1193–1205.
- 1194 Mullainathan, S. & Shleifer, A. (2005). The market for news. *American Eco-*  
1195 *nomics Review*, **95**, 1031–1053.
- 1196 Newman, N., Fletcher, R., Kalogeropoulos, A., Levy, D.A.L. & Nielsen, R.  
1197 (2018). *Reuters institute digital news report 2018*, vol. 2018. Reuters Insti-  
1198 tute for the Study of Journalism.
- 1199 Newman, N., Fletcher, R., Kalogeropoulos, A. & Nielsen, R. (2019). *Reuters*  
1200 *institute digital news report 2019*, vol. 2019. Reuters Institute for the Study  
1201 of Journalism.



- 1202 Noble, S.U. (2018). *Algorithms of Oppression: How search engines reinforce*  
1203 *racism*. NYU Press.
- 1204 Otterbacher, J., Bates, J. & Clough, P. (2017). Competent men and warm  
1205 women: Gender stereotypes and backlash in image search results. *Proceed-*  
1206 *ings of the 2017 chi conference on human factors in computing systems*, pp.  
1207 6620–6631.
- 1208 Otterbacher, J., Checco, A., Demartini, G. & Clough, P. (2018). Investigating  
1209 user perception of gender bias in image search: the role of sexism. *The*  
1210 *41st International ACM SIGIR Conference on Research & Development in*  
1211 *Information Retrieval*, pp. 933–936.
- 1212 Pan, B., Hembrooke, H., Joachims, T., Lorigo, L., Gay, G. & Granka, L. (2007).  
1213 In google we trust: Users’ decisions on rank, position, and relevance. *Journal*  
1214 *of computer-mediated communication*, **12**, 801–823.
- 1215 Rábiger, S., Gezici, G., Saygın, Y. & Spiliopoulou, M. (2018). Predicting  
1216 worker disagreement for more effective crowd labeling. *2018 IEEE 5th In-*  
1217 *ternational Conference on Data Science and Advanced Analytics (DSAA)*,  
1218 pp. 179–188, IEEE.
- 1219 Raji, I.D. & Buolamwini, J. (2019). Actionable auditing: Investigating the im-  
1220 pact of publicly naming biased performance results of commercial ai prod-  
1221 ucts. *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and*  
1222 *Society*, pp. 429–435.
- 1223 Robertson, R.E., Jiang, S., Joseph, K., Friedland, L., Lazer, D. & Wilson, C.  
1224 (2018a). Auditing partisan audience bias within google search. *Proceedings*  
1225 *of the ACM on Human-Computer Interaction*, **2**, 148.
- 1226 Robertson, R.E., Lazer, D. & Wilson, C. (2018b). Auditing the personalization  
1227 and composition of politically-related search engine results pages. *Proceed-*  
1228 *ings of the 2018 World Wide Web Conference*, pp. 955–965.
- 1229 Saez-Trumper, D., Castillo, C. & Lalmas, M. (2013). Social media news com-  
1230 munities: gatekeeping, coverage, and statement bias. *Proceedings of the 22nd*  
1231 *ACM international conference on Information & Knowledge Management*,  
1232 pp. 1679–1684, ACM.
- 1233 Sandvig, C., Hamilton, K., Karahalios, K. & Langbort, C. (2014). Auditing  
1234 algorithms: Research methods for detecting discrimination on internet plat-  
1235 forms. *Data and discrimination: converting critical concerns into productive*  
1236 *inquiry*, **22**.
- 1237 Sapiezynski, P., Zeng, W., E Robertson, R., Mislove, A. & Wilson, C. (2019).  
1238 Quantifying the impact of user attention on fair group representation in  
1239 ranked lists. *Companion Proceedings of The 2019 World Wide Web Con-*  
1240 *ference*, pp. 553–562.
- 1241 Sarcona, C. (2019). Organic search click through rates:  
1242 The numbers never lie. [https://www.zerolimitweb.com/  
1243 organic-vs-ppc-2019-ctr-results-best-practices/](https://www.zerolimitweb.com/organic-vs-ppc-2019-ctr-results-best-practices/), accessed: 2019-  
1244 09-06.
- 1245 Stokes, P. (2019). False media balance. [https://www.newphilosopher.com/  
1246 articles/false-media-balance/](https://www.newphilosopher.com/articles/false-media-balance/), accessed: 2019-09-15.

- 1247 Su, H., Deng, J. & Fei-Fei, L. (2012). Crowdsourcing annotations for visual  
1248 object detection. *Workshops at the Twenty-Sixth AAAI Conference on Ar-*  
1249 *tificial Intelligence*.
- 1250 Tavani, H. (2012). Search engines and ethics, .
- 1251 Vincent, N., Johnson, I., Sheehan, P. & Hecht, B. (2019). Measuring the im-  
1252 portance of user-generated content to search engines. *Proceedings of the*  
1253 *International AAAI Conference on Web and Social Media*, vol. 13, pp. 505–  
1254 516.
- 1255 Vondrick, C., Patterson, D. & Ramanan, D. (2013). Efficiently scaling up  
1256 crowdsourced video annotation. *International Journal of Computer Vision*,  
1257 **101**, 184–204.
- 1258 White, R. (2013). Beliefs and biases in web search. *Proceedings of the 36th*  
1259 *international ACM SIGIR conference on Research and development in in-*  
1260 *formation retrieval*, pp. 3–12, ACM.
- 1261 Yang, K. & Stoyanovich, J. (2017). Measuring fairness in ranked outputs.  
1262 *Proceedings of the 29th International Conference on Scientific and Statistical*  
1263 *Database Management*, p. 22, ACM.
- 1264 Yigit-Sert, S., Altıngövdü, I.S. & Ulusoy, Ö. (2016). Towards detecting media  
1265 bias by utilizing user comments, .
- 1266 Yuen, M.C., King, I. & Leung, K.S. (2011). A survey of crowdsourcing systems.  
1267 *2011 IEEE Third International Conference on Privacy, Security, Risk and*  
1268 *Trust and 2011 IEEE Third International Conference on Social Computing*,  
1269 pp. 766–773, IEEE.
- 1270 Zehlike, M., Bonchi, F., Castillo, C., Hajian, S., Megahed, M. & Baeza-Yates,  
1271 R. (2017). Fa\* ir: A fair top-k ranking algorithm. *Proceedings of the 2017*  
1272 *ACM on Conference on Information and Knowledge Management*, pp. 1569–  
1273 1578, ACM.