Towards Fair RAG: On the Impact of Fair Rank-ING IN RETRIEVAL-AUGMENTED GENERATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Many language models now enhance their responses with retrieval capabilities, leading to the widespread adoption of retrieval-augmented generation (RAG) systems. However, despite retrieval being a core component of RAG, much of the research in this area overlooks the extensive body of work on fair ranking, neglecting the importance of considering all stakeholders involved. This paper presents the first systematic evaluation of RAG systems integrated with fair rankings. We focus specifically on measuring the fair exposure of each relevant item across the rankings utilized by RAG systems (i.e., item-side fairness), aiming to promote equitable growth for relevant item providers. To gain a deep understanding of the relationship between item-fairness, ranking quality, and generation quality in the context of RAG, we analyze nine different RAG systems that incorporate fair rankings across seven distinct datasets. Our findings indicate that RAG systems with fair rankings can maintain a high level of generation quality and, in many cases, even outperform traditional RAG systems, despite the general trend of a tradeoff between ensuring fairness and maintaining system-effectiveness. We believe our insights lay the groundwork for responsible and equitable RAG systems and open new avenues for future research. We publicly release our codebase and dataset.

1 INTRODUCTION

030 031

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

029

032 In recent years, the concept of fair ranking has emerged as a critical concern in modern information 033 access systems (Ekstrand et al., 2022). However, despite its significance, fair ranking has yet to be 034 thoroughly examined in the context of retrieval-augmented generation (RAG) (Lewis et al., 2020; Asai et al., 2024), a rapidly advancing trend in natural language processing (NLP) systems (Kim 035 et al., 2024). To understand why this is important, consider the RAG system in Figure 1, where a 036 user asks a question about running shoes. A classic retrieval system might return several documents 037 containing information from various running shoe companies. If the RAG system only selects the top two documents, then information from the remaining two relevant companies will not be relayed to the predictive model and will likely be omitted from its answer. The fair ranking literature refers 040 to this situation as unfair because some relevant companies (i.e., in documents at position 3 and 4) 041 receive less or no exposure compared to equally relevant company in the top position (Ekstrand et al., 042 2022). 043

Understanding the effect of fair ranking in RAG is fundamental to ensuring responsible and equitable 044 NLP systems. Since retrieval results in RAG often underlie response attribution (Gao et al., 2023), unfair exposure of content to the RAG system can result in incomplete evidence in responses (thus 046 compromising recall of potentially relevant information for users) or downstream representational 047 harms (thus creating or reinforcing biases across the set of relevant entities). In situations where 048 content providers are compensated for contributions to inference, there can be financial implications for the unfairness (Balan et al., 2023; Lyu et al., 2023; Henderson et al., 2023). Indeed, the fair ranking literature indicates that these are precisely the harms that emerge when *people* are searchers 051 (Ekstrand et al., 2022), much less RAG systems, where the searchers are machines. RAG complicates these challenges since it often truncates rankings to much shorter lengths to fit the generator's limited 052 context size (Bahri et al., 2020; Hofstätter et al., 2023; Kim et al., 2024), making equal exposure of relevant items even harder.



063

064 Figure 1: Fairness concerns in RAG. A simplified example of how RAG models can ignore equally 065 relevant items (d_3 and d_4) and always consume the fixed top-scoring items (d_1 and d_2) with the same 066 order of ranking over the multiple user requests. This is due to the deterministic nature of the retrieval 067 process and a short context-length of a language model that necessitates the top-k truncation of a 068 ranked list.

069

Moreover, the fact that machines are now the searchers necessitates a different notion of itemworthiness (how deserving an item is to be included in a ranked list). Traditionally, ranking quality 071 has been assessed based on relevance labels, which are created according to how relevant an item is to 072 the user's query (Saracevic, 2016). However, with RAG systems, where the consumer is a language 073 model, there is a growing shift towards evaluating ranking quality based on utility labels, which 074 are determined by the usefulness of an item in aiding the model's task performance, rather than its 075 relevance to the query (Salemi & Zamani, 2024a; Zhang et al., 2024a). 076

077 This shift from relevance to utility in the concept of item-worthiness can significantly alter our understanding of the relationship between fairness and ranking quality (Balagopalan et al., 2023) particularly the tradeoffs that are well-known in the fair ranking literature (Biega et al., 2018; Diaz 079 et al., 2020; Singh & Joachims, 2019). Since previous fair ranking studies were conducted based on relevance judgments, they may need to be reexamined in light of utility-based judgments within the 081 context of RAG.

083 Our research aims to bridge the gap between traditional fair ranking studies and the emerging changes posed by RAG systems, ultimately enhancing our understanding of the interplay between fairness, 084 ranking quality, and the effectiveness of RAG systems. We do this by evaluating RAG systems with a 085 fairness-aware retriever across seven different tasks, experimenting with varying levels of retrieval fairness to observe changes in ranking quality and generation quality (utility).¹ 087

Our empirical results show that, in the context of machine users, there also exists an overall trend of fairness-quality tradeoff with respect to both retrieval and generation quality. However, the magnitude of this tradeoff is not particularly severe. In fact, we find that RAG models equipped with a fair 090 ranker can often preserve a significant level of retrieval and generation quality, and in some cases, 091 even surpass the quality achieved by the traditional RAG setup with a deterministic ranker that lacks 092 fairness considerations.

094 This surprising finding offers significant insight into the potential of RAG-based applications, suggesting that fair treatment of individual content providers can be achieved without sacrificing much 095 of the high-quality service delivered to end-users. This challenges the conventional assumption of an 096 inevitable tradeoff between fairness and quality, opening new avenues for developing more equitable and effective RAG systems. 098

099 100

101

2 **BACKGROUND & RELATED WORK**

102 **Retrieval-Augmented Generation**. RAG, a specific type of retrieval-enhanced machine learning 103 (REML) (Zamani et al., 2022; Kim et al., 2024), has been widely adopted in various domains, 104 including language modeling (Khandelwal et al., 2020), question-answering (Izacard et al., 2023), 105 personalization (Salemi et al., 2024b;a; Kumar et al., 2024; Neelakanteswara et al., 2024), and

¹Throughout this paper, we use "utility" and "generation quality" interchangeably to refer to the downstream effectiveness of RAG models, measured by arbitrary string utility metrics.



119 Figure 2: Experimental design to investigate the impact of item-fairness on ranking and generation 120 quality in RAG. To evaluate system performance across multiple identical user requests, we sample 121 N rankings from a stochastic retriever. We then measure the fairness and quality of the sampled 122 rankings (Expected Exposure) and assess the system's expected end-performance (Expected Utility). 123 The query and prompt generators are omitted in the figure for brevity. Details on implementing a 124 RAG system with fair rankings in a production environment can be found in Appendix E.

125 126

111

113

114

recommendation (Zeng et al., 2024). Studies on the evaluation of RAG models have primarily 127 focused on their effectiveness, including end-to-end performance (Izacard et al., 2023; Guu et al., 128 2020; Lewis et al., 2020) and the assessment of individual components (Es et al., 2024; Saad-Falcon 129 et al., 2024; Salemi & Zamani, 2024a). However, little research has focused on evaluating fairness 130 in retrieval-enhanced generation models, with the exception of recent work (Shrestha et al., 2024), 131 which improved demographic diversity in human image generation by conditioning a generative 132 model with externally retrieved images that help debias the generation process. 133

Fairness in Ranking. Fair ranking has been approached through various definitions based on 134 normative concerns, primarily with distinctions made according to the stakeholders we prioritize. 135 These include consumer-side fairness (Mehrotra et al., 2017; Ekstrand et al., 2024), which focuses on 136 how fairly a system delivers satisfaction to users; provider-side fairness (Sapiezynski et al., 2019; 137 Jaenich et al., 2024), which addresses how fairly item providers receive monetary or reputational 138 rewards; and item-side fairness (Jiang et al., 2024), which examines how fairly items are treated in 139 terms of representation or exposure. The motivation of item-side fairness is closely linked to provider-140 side fairness, as unfair treatment of items can lead to unfair compensation for providers. These 141 fairness concerns can be further categorized by the scope of stakeholders, encompassing individual 142 fairness—ensuring similar treatment for similar individuals—and group fairness—ensuring equitable outcomes across different groups (Caton & Haas, 2020; Ekstrand et al., 2022). Previous studies 143 have focused on developing metrics to measure fairness (Raj & Ekstrand, 2020) and optimizing 144 fair retrievers within a single (Yang & Stoyanovich, 2017; Zehlike et al., 2017; Sapiezynski et al., 145 2019) or multiple rankings (Diaz et al., 2020; Singh & Joachims, 2018; Biega et al., 2018; Singh 146 & Joachims, 2019). In the context of provider- and item-side fairness, ensuring equal exposure 147 of similar items across multiple rankings has gained significant attention (Ekstrand et al., 2022). 148 To achieve this, researchers have used stochastic rankers that return a distribution of rankings, in 149 contrast to deterministic rankers commonly found in areas like RAG, which produce a fixed ranking. 150 This approach ensures that, in expectation, similar items receive equal exposure across multiple 151 user requests, with the distributions typically based on the merit of the rankings, such as an item's 152 relevance (Diaz et al., 2020; Singh & Joachims, 2019).

153 In this research, we employ a stochastic ranker in RAG to enhance *individual item-side fairness*, 154 aiming to ensure equal expected exposure for items that offer similar merits.

155 156

3 EXPERIMENTAL METHODOLOGY

157 158

In traditional RAG systems, a user input is used to query a retrieval system for recommended items 159 from some corpus, which are then used for generation. Given user input x, a query q generated by the query generation function $\phi_q(x)$, and a corpus of documents C, a *deterministic retriever* 161 $\mathcal{R}(q, \mathbb{C})$ returns a fixed ranked list L every time q is seen. Retrieval is followed by a top k truncation which is passed to a prompt generation function $\phi_p(x, L_{1:k})$ that returns a final prompt \bar{x} , which is subsequently passed to the language model $\mathcal{G}(\bar{x})$. Because deterministic retrievers allocate exposure to the same item over repeated samples, RAG systems with deterministic retrievers present a challenge to ensuring equal exposure of relevant items to the generator.

166 To address the issue of unfairness in the rankings passed to the generator, we can convert a determin-167 istic retriever into a stochastic retriever, which can, in expectation, provide fair rankings (Diaz et al., 168 2020). By sampling a ranking based on its quality to users—in this case generators—the expected 169 exposure of different relevant items becomes similar and, therefore, fairer (Appendix E). Because 170 decisions are stochastic, the fairness and quality of stochastic retrieval is evaluated based on a sample 171 of rankings. Similarly, since a sampled ranking is processed by a generator, we also compute the 172 expected generator effectiveness over sampled rankings. The complete evaluation pipeline of a RAG system with a stochastic retriever is illustrated in Figure 2. 173

The following sections describe how we construct a test collection with utility labels (§3.1), how we stochastically sample multiple rankings (§3.2), and how we evaluate the fairness and ranking quality of the sampled rankings (§3.3.1) and measure the effectiveness of a RAG system given multiple rankings (§3.3.2).

179 3.1 CONSTRUCTION OF A TEST COLLECTION WITH UTILITY LABELS

Setting an appropriate proxy for measuring item-worthiness is crucial in the evaluation of fairness
(Balagopalan et al., 2023). Drawing on the insight that utility-based judgments are more suitable
than relevance judgments in the context of RAG (Zhang et al., 2024a; Salemi & Zamani, 2024a), we
annotate item-level utility labels for all items in the corpus.

185 We define an item's worthiness by the marginal gain in utility (utility-gain) it provides to a language model (specifically, the generator in a RAG system) when used to solve a specific task as part of the augmentation process. To assess this utility-gain, each item in the corpus is individually supplied 187 to the generator along with an input question. The utility-gain is then calculated as the difference 188 between the utility of the augmented generator and that of a baseline language model without any 189 information about the item. Formally, let u_i denote the baseline string utility score from the vanilla 190 language model prompted only with the input question, and let u_i represent the utility score from the 191 language model with a prompt augmented by the j'th item d_i in the corpus. The item d_i is considered 192 useful if the utility-gain $\delta_i = u_i - u_i$ is positive, and not useful otherwise (see Appendix B). 193

Therefore, the item-level utility labels are designed to be both task- and generator-dependent, as the utility of each item varies depending on the task and the language model used. This labeling process also aligns with the principles of task-based information retrieval, where, in the context of *human* searchers, document utility may vary on how the user expects to use the document (Kelly et al., 2013).

198 199

178

180

3.2 FAIRNESS-AWARE STOCHASTIC RETRIEVER

200 Stochastic retrievers have been used for various purposes, such as optimization of retrieval models 201 (Bruch et al., 2020; Zamani & Bendersky, 2024; Guiver & Snelson, 2009; Oosterhuis, 2021), as well 202 as ensuring equitable exposure of items (Oosterhuis, 2022; Diaz et al., 2020; Oosterhuis, 2021). Many 203 of these studies use Plackett-Luce sampling (Plackett, 1975) to achieve the stochasticity of retrieval. 204 We follow the line of research and formally define how we derive a fairness-aware stochastic retriever through Plackett-Luce sampling. To enhance sampling efficiency, we adopt the methodology of 205 Oosterhuis (2021), and for controllable randomization, we utilize the approach proposed by Diaz 206 et al. (2020). 207

Given *n* items in a corpus C, a vector of retrieval scores $s \in \mathbb{R}^n$ can be obtained from $\mathcal{R}(q, C)$, which can be used to generate a ranked list *L*. We then min-max normalize retrieval scores to be in [0, 1] in order to construct a multinomial distribution over items (Biega et al., 2018). The probability of an item *d* being selected as the *i*'th item in a new ranking π through Plackett-Luce sampling is given by

214

$$p(d|L_{1:i-1}) = \frac{\exp(\mathbf{s}_d) \mathbb{I}[d \notin L_{1:i-1}]}{\sum_{d' \in \mathbb{C} \setminus L_{1:i-1}} \exp(\bar{\mathbf{s}}_{d'})}$$
(1)

 $(-) = \begin{bmatrix} 1 \\ - \end{bmatrix} + T$

where $L_{1:i-1}$ is the partial ranking up to position i - 1, \bar{s} represents the normalized retrieval score vector, and \bar{s}_d is the normalized score of item d. Using this probability, we iteratively sample an

item, set its probability to 0, renormalize the distribution, and repeat the process. The probability of generating a complete ranking is then given by the product of the placement probabilities for each item, i.e., $p(\pi|q) = \prod_{i=1}^{n} p(\pi_i|\pi_{1:i-1})$.

This repeated sampling and renormalization process can be efficiently managed using the Gumbel-Softmax trick (Gumbel, 1954; Maddison et al., 2017), which enables the sampling of rankings to be performed at the speed of sorting (Oosterhuis, 2021). To do so, for each sampling iteration, we draw $U_i \sim \text{Uniform}(0, 1)$, followed by generating a Gumbel noise $G_i = -\log(-\log(U_i))$. The probability of each sampled ranking is then obtained by sorting the items based on their perturbed scores $\tilde{\mathbf{s}}_{d_i} = \bar{\mathbf{s}}_{d_i} + G_i$.

226

227

234 235 236

245

3.2.1 CONTROLLING THE LEVEL OF FAIRNESS

Adjusting the level of randomization directly controls the degree of item-fairness, aligning with our goal to observe how varying levels of fairness in rankings affect the ranking and generation quality of a RAG model. To obtain the controllability, we follow the work of Diaz et al. (2020) and use a temperature parameter α . We apply the scalar α to each value in the normalized score vector \bar{s} by raising each value to the power of α .² This process is done before the scores are passed to the sampling policy. Therefore, the modified sampling distribution is thus defined as:

$$p(d|L_{1:i-1}) = \frac{\exp(\bar{\mathbf{s}}_d^{\alpha})\mathbb{1}[d \notin L_{1:i-1}]}{\sum_{d' \in \mathbb{C} \setminus L_{1:i-1}} \exp(\bar{\mathbf{s}}_{d'}^{\alpha})}$$
(2)

237 This implies that the sharpness of the sampling distribution is controlled by the α . A higher α 238 amplifies the probability of items with higher retrieval scores being sampled. Therefore, if multiple 239 rankings are sampled by the stochastic retriever with high α , it results in high disparity (i.e., item-side 240 unfairness) of sampled rankings. At extreme, with considerably high α , the procedure results in the 241 identical rankings which is the behavior of a deterministic ranker (i.e., maximum item-unfairness). On the other hand, a lower α reduces the disparity of sampled rankings, making the exposure distribution 242 fairer. At extreme, when $\alpha = 0$, the sampling procedure becomes uniformly random and achieves the 243 lowest disparity (i.e., maximum item-fairness) in the sampled rankings. 244

246 3.3 EVALUATION

247 As mentioned in Section 3, because we are dealing with stochastic retrievers, we need to measure 248 the *expected* behavior of the system. Let $\mathcal{S}(\mathbf{s}, N, k)$ be the stochastic sampler that samples a set of 249 N rankings $\sigma = \{\pi\}$, given the retrieval scores s, where each ranking π is truncated to the size of 250 k. From each ranking, we can get an output $\hat{y} = \mathcal{G}(\phi_p(x,\pi))$. With an arbitrary fairness metric 251 $\mu_f(\sigma)$ and a ranking quality metric $\mu_r(\sigma)$ that takes a set of rankings as an input, we can measure the 252 degree of fairness and ranking quality of the sampled rankings. Similarly, an arbitrary string utility 253 metric $\mu_u(y, \hat{y})$, such as ROUGE, can be used to assess an expected effectiveness of a RAG system 254 by calculating the average of the N metric scores.

In this paper, based on the empirical investigation done by Raj & Ekstrand (2020), we use expected exposure disparity (EE-D) and expected exposure relevance (EE-R) (Diaz et al., 2020) as μ_f and μ_r , respectively (§3.3.1). For μ_u , we select the metric depending on the task, and we get the expectation of the utility of a RAG model which we call an expected utility (EU) (§3.3.2).

260 261 3.3.1 EXPECTED EXPOSURE IN THE CONTEXT OF MACHINE USERS

Expected Exposure (EE) (Diaz et al., 2020) works by estimating the exposure of items across rankings (e.g., σ) created by a subject model, and comparing them with an optimal set of rankings that always satisfy the item-fairness. To represent the attention over *n* items given by the consumer (generator in RAG), an $n \times 1$ system exposure vector ϵ is created. This is then compared with an $n \times 1$ target exposure vector ϵ^* , where it represents the exposure of items allocated by an oracle retriever that always rank useful items above non-useful ones (Diaz et al., 2020).

²We normalized the values to the range of [1, 2] instead of [0, 1]. The addition of 1 effectively serves the same purpose as adjusting a real-numbered α . We chose this range to allow for an integer-valued α .

With the system and target exposure vector $\epsilon \in \mathbb{R}^n$ and $\epsilon^* \in \mathbb{R}^n$, we can get the difference between the two by the squared l2 distance:

$$\|\epsilon - \epsilon^*\|_2^2 = \|\epsilon\|_2^2 - 2\langle\epsilon, \epsilon^*\rangle + \|\epsilon^*\|_2^2$$
(3)

273 274

This difference yields two metrics useful for fairness and ranking quality evaluation. $\|\epsilon\|_2^2$ can be a measure for disparity of rankings (EE-D), and $\langle \epsilon, \epsilon^* \rangle$ can be a measure of ranking quality (EE-R) by calculating the degree of alignment of system exposure to the target exposure (i.e., how much of the exposure is on useful items). Therefore, the higher the value of EE-D, the more unfair the set of rankings are, and the higher the value of EE-R, the closer the set of system rankings are to the optimal set of rankings with respect to the ranking quality.

280 The exposure of an item is calculated by modeling users' (e.g., generators in RAG) attention to each 281 item in a ranking. For example, one can assume that the user is affected by position bias and gives 282 attention following an exponential decay (Moffat & Zobel, 2008). However, these browsing models 283 were developed for human-users not for machine-users, so we need a different user behavior model 284 for generators in RAG. For simplicity, we assume that the machine-user can consume the items by 285 giving equal attention to all the items that were passed to the context, but pays 0 attention to the items placed after the k'th position due to the top-k truncation. This makes the user browsing model a step 286 function parameterized by k. In this work, a relevance-independent machine-user model (MU) is set 287 to the step function that reflects the behavior of *top-k* truncation of RAG: 288

$$MU(i) = \begin{cases} 1 & \text{if } i \le k \\ 0 & \text{otherwise} \end{cases} = \mathbb{1}[i \le k]$$
(4)

Given this machine user browsing model and a mapping from item index to its rank denoted as $\bar{\pi}_d$, a system exposure for each item d is calculated as

$$\epsilon_d = \sum_{\pi \in S_n} p(\pi|q) \mathrm{MU}(\bar{\pi}_d) \tag{5}$$

and target exposures for a useful item d and a unuseful item d^- are calculated as

N

$$\epsilon_d^* = \frac{1}{m} \sum_{i=1}^m \mathrm{MU}(i) = \begin{cases} 1 & \text{if } m \le k \\ \frac{k}{m} & \text{otherwise} \end{cases} \qquad \epsilon_{d^-}^* = \begin{cases} \frac{k-m}{n-m} & \text{if } m \le k \\ 0 & \text{otherwise} \end{cases}$$
(6)

299 300 301 302

309 310

298

289

290 291

295 296 297

3.3.2 EXPECTED UTILITY

Given the set of N sampled rankings σ , we individually augment the generator with each ranking $\pi \in \sigma$, resulting in N outputs from the generator. The utility of these outputs is then measured using an arbitrary string utility metric μ_u . To determine the anticipated utility of a RAG model with fair rankings—represented by the tuple of a stochastic ranking sampler S and a generator \mathcal{G} —we calculate the expected utility (EU) of the RAG system given an instance x.

$$\operatorname{EU}(\langle \mathcal{S}, \mathcal{G} \rangle | x) = \mathbb{E}_{\pi \sim \mathcal{S}}[\mu_u(y, \hat{y}_\pi)] = \sum_{\pi \in S_n} p(\pi | q) \mu_u(y, \hat{y}_\pi) \approx \frac{1}{N} \sum_{\pi \in \sigma} \mu_u(y, \hat{y}_\pi)$$
(7)

where \hat{y}_{π} is the prediction of a system given the ranking π , S_n is the symmetric group of a ranked list *L* from the deterministic retriever \mathcal{R} , and $\sum_{\pi \in S_n} p(\pi | q) = 1$.

314 3.3.3 NORMALIZATION OF METRICS315

From Equation 3, we decompose the metric into EE-D and EE-R. Since the bounds of these metrics depend on the number of useful items, normalization must be applied per query. Both metrics are min-max normalized based on their theoretical lower and upper bounds. We denote the normalized EE-D and EE-R as $\overline{\text{EE-D}}_q$ and $\overline{\text{EE-R}}_q$, respectively.

However, theoretically determining the bounds of the expected utility (EU) of a RAG model is
challenging. To address this, we normalized the EU by the model's empirical upper bound, the
maximum observed utility across all runs of the experiment with the same generators. To approach
the true upper bound, these runs include RAG models with an oracle retriever that consistently ranks
useful items (i.e., those with positive utility labels) above non-useful ones, stochastically returning

one of the m!(n-m)! different rankings, where m represents the number of useful items in the corpus. We denote the normalized EU as \overline{EU}_q , which can be interpreted as the distance to the optimal utility. From this section onwards, we omit the symbol q from the normalized metrics for brevity. Proofs and details on how each metric is normalized by its lower and upper bounds can be found in the Appendix. C.

330 331

4 EXPERIMENT SETUP

332 We choose the LaMP benchmark (Salemi et al., 2024b) for our dataset. It assesses the personalization 333 capability of language models through retrieval-augmentation of users' interaction history in a 334 platform. LaMP includes various prediction tasks, such as classification, regression, and generation, 335 and is well-suited for tasks where multiple items can be relevant/useful, unlike QA tasks with typically 336 one or two provenance items. The retrieval items in LaMP have clear providers and consumers, 337 aligning with our goal to ensure fairness for individual item providers. For example, in LaMP-1, retrieval items are academic papers, where exposure can increase citation counts for authors. In 338 LaMP-4, retrieval items are news articles, where exposure can lead to monetary compensation for 339 journalists. Due to the absence of a test set, we constructed a test collection as described in §3.1, 340 using the first thousand entries of a user-based development set. Then, we discarded entries that have 341 only one useful item in the corpus, as it is unnecessary to concern item-fairness in that case. We 342 release the test collection, and the dataset statistics can be found in the Appendix J. 343

We use BM25 (lexical retriever) (Robertson et al., 1995), SPLADE (learned sparse retriever), and Contriever (biencoder dense retriever) (Izacard et al., 2022) as deterministic retrievers providing retrieval scores to base the sampling on. These models represent commonly used retrievers in the RAG literature (Kim et al., 2024). We use a sampling size of N = 100 and a truncation size of k = 5.

For generation models, we use Flan-T5-Small, Flan-T5-351 Base, and Flan-T5-XXL (Chung et al., 2022). For de-352 coding strategy, beam size is set to 4, and no sampling 353 strategy is used. This is to ensure that stochasticity is only 354 introduced to the retriever for controlled experiments. Full 355 implementation details can be found in Appendix K. With 356 the three base retrievers and three generators, we configure 357 nine different RAG models and evaluate them on the seven 358 LaMP tasks. Utility measurement of the generated strings 359 followed the metrics used in the LaMP paper.



Figure 3: Effect of a temperature parameter α on the disparity of rankings, in the LaMP task 4, a text generation task. The RAG model is configured with the Contriever and Flan-T5-Base. Each data point represents the normalized EE-D of each run of the experiment (i.e., one query ->N sampled rankings -> EE-D of the N rankings).

We repeat the experiments with four different temperature parameters $\alpha = 1, 2, 4, 8$, which allows us to assess the utility of the RAG models with different levels of item-

fairness. From Figure 3, we observe how effectively α , described in the Equation 2, controls the disparity of rankings. For example, when α is set to 4, we usually obtain a set of sampled rankings with $\overline{\text{EE-D}}$ mostly in the range of [0.5, 0.8], and when α is set to 8, we often get a set of sampled rankings with $\overline{\text{EE-D}} = 1$. Refer to Appendix D to see the full description of the effect of α on the other metrics.

369 5 R

370 371

368

5 Results

RQ1: Is there a tradeoff between ensuring item-fairness in rankings and maintaining high ranking
 quality when utility labels are used for evaluation?

By gathering all four repeated runs of the experiments with different α values, we can plot the trend of ranking quality ($\overline{\text{EE-R}}$) against item fairness ($\overline{\text{EE-D}}$), as shown in Figure 4.

As shown in previous studies (Singh & Joachims, 2019; Diaz et al., 2020), there is a well-known
 tradeoff between fairness and ranking quality for human users. Similarly, we observe a general
 tradeoff for machine users. However, unlike past findings, this tradeoff is not always strict. For

instance, in Figure 4, both SPLADE and Contriever maintain consistently high ranking quality while
 being considerably fairer, and for BM25, ranking quality even improves as fairness increases, up to a
 certain point.

At the rightmost side of the lines, where $\overline{\text{EE-D}} = 1$ (rep-382 resenting the performance of deterministic rankers), we 383 observe that these rankers do not always deliver the high-384 est ranking quality. This suggests that commonly used 385 deterministic rankers in RAG systems may be suboptimal, 386 and that ranking quality can be improved while ensuring 387 item fairness. This becomes even clearer when examining 388 the impact of fair ranking on the downstream performance of a RAG system. 389

The leftmost side of the lines, where $\overline{\text{EE-D}} = 0$, represents the performance of a uniformly random ranking policy. At this point, the measured ranking quality should approximate the proportion of positively labeled items in the corpus, which is 31% based on data statistics (Appendix



Figure 4: Relationship between itemfairness and the quality of rankings.

J). This is notably higher than in non-RAG (human-user) settings, where the percentage of relevant documents is typically much smaller, resulting in a $\overline{\text{EE-R}}$ value near 0 (Diaz et al., 2020).

To quantify the performance of fair rankers, we calculate the area under the disparity-ranking quality curve (Figure 4), with higher values indicating stronger ranking quality. We also measure the tradeoff by fitting a linear line to the experiment results, where a steeper slope reflects a stronger tradeoff between fairness and ranking quality. Based on these metrics, we observe that Contriever-based models exhibit the highest tradeoff, while BM25-based models show the lowest, despite their poor retrieval quality. Overall, SPLADE-based models achieve high retrieval quality while maintaining a relatively low tradeoff. For detailed plots and quantifications, refer to Appendix F.



Figure 5: (a) Strong correlation between ranking quality and generation quality. (b) Fairnessgeneration quality tradeoff. Full plots and quantifications are provided in Appendix G and H.

Fairness Intervals				Fairness Intervals							
Model (baseline utility)	[0.0, 0.2)	[0.2, 0.4) [0.4, 0.6) [0.6, 0.8) [0.8, 1.0) Mode		Model (baseline utility)	[0.0, 0.2)	[0.2, 0.4)	[0.4, 0.6)	[0.6, 0.8)	[0.8, 1.0)		
LaMP-1						LaN	/IP-4				
BM25+FlanT5Small (0.308)	-0.12	-0.13	-0.18	-0.02	-0.15	BM25+FlanT5Small (0.217)	-0.06	0.00	+0.02	+0.01	0.00
BM25+FlanT5Base (0.670)	-0.20	-0.04	-0.08	-0.05	-0.02	BM25+FlanT5Base (0.223)	-0.06	0.00	+0.03	+0.01	+0.02
BM25+FlanT5XXL (0.531)	-0.07	+0.03	+0.02	+0.06	+0.11	BM25+FlanT5XXL (0.322)	-0.05	+0.11	+0.03	+0.03	+0.05
SPLADE+FlanT5Small (0.241)	-0.03	-0.22	+0.19	-0.04	+0.14	SPLADE+FlanT5Small (0.235)	-0.07	-0.01	+0.02	+0.03	+0.02
SPLADE+FlanT5Base (0.646)	-0.15	+0.06	+0.08	0.00	+0.03	SPLADE+FlanT5Base (0.268)	-0.10	-0.03	+0.02	0.00	+0.02
SPLADE+FlanT5XXL (0.671)	-0.18	-0.16	+0.05	+0.02	+0.01	SPLADE+FlanT5XXL (0.342)	-0.06	+0.09	+0.05	+0.03	+0.04
Contriever+FlanT5Small (0.286)	-0.08	-0.29	-0.06	+0.03	-0.14	Contriever+FlanT5Small (0.254)	-0.09	-0.02	0.00	+0.01	0.00
Contriever+FlanT5Base (0.637)	-0.16	+0.05	-0.06	+0.03	0.00	Contriever+FlanT5Base (0.268)	-0.10	-0.02	+0.01	0.00	+0.01
Contriever+FlanT5XXL (0.651)	-0.19	-0.04	-0.11	+0.03	0.00	Contriever+FlanT5XXL (0.367)	-0.09	+0.06	+0.01	+0.01	+0.03

428 429

419

420

Table 1: Each value in the table is the difference between the utility of a baseline (deterministic)
 RAG model and the average utility of a fairer RAG model at a specific fairness interval. Nonnegative differences are highlighted. Full results are listed in Appendix I.

RQ2: Is there a tradeoff between ensuring item-fairness in ranking and maintaining high generation quality of a RAG model?

Before examining the relationship between fairness and RAG utility, Figure 5a shows an auxiliary result confirming a strong correlation between utility-based ranking quality and the effectiveness of RAG models. This is unsurprising, as item-worthiness judgments were based on the utility-gain provided by the generator. However, this correlation suggests that the tradeoff observed in the disparity-ranking quality curve (Figure 4) is likely to manifest similarly due to this strong relationship.

In fact, as observed from the disparity-utility curve (Figure 5b), we see a global trend of a non-strict tradeoff (i.e., RAG models maintain high generation quality while being considerably fair, and often even achieve higher quality).

However, a closer look at the local trend offers a significant insight: *RAG systems with fair ranking can* often achieve higher system-effectiveness compared to models with deterministic rankers. In Table 1, we divided the fairness levels into five intervals based on the normalized EE-D. As shown in the table and Appendix I, improving fairness to the level of $\overline{\text{EE-D}} \in [0.8, 1.0)$, and even $\overline{\text{EE-D}} \in [0.6, 0.8)$, can often enhance the utility of many RAG models across most LaMP tasks. For example, having $\overline{\text{EE-D}}$ in the range of [0.8, 1.0) outperforms the baseline for all models in LaMP-2 and for seven out of nine models in LaMP tasks 4, 5, and 6.

449 450

451 452

6 DISCUSSION AND CONCLUSION

Why do we often see higher system utility in fairer RAG models? Although there is a general trend of a fairness-utility tradeoff, we observe that certain levels of fairness can actually improve the utility of a baseline RAG model. Recent line of research have uncovered relevant findings: 1) generators are not robust to changes in the position of useful information (Liu et al., 2024); 2) items with high retrieval scores often include distracting content that can reduce the system-effectiveness (Cuconasu et al., 2024; Ru et al., 2024); and 3) introducing some random documents can significantly boost the utility of RAG (Cuconasu et al., 2024).

Building on these existing results, we find that perturbing the initial ranking through stochastic sampling often can impact the performance of certain inference decisions and lead to changes in the system's expected end-performance. In our experiments, we observe that the expected utility generally increases within the fairness interval of [0.8, 1.0). This suggests that a fixed ranking from a deterministic ranker may be suboptimal for the generator, and that perturbing the ranking, along with the repositioning of items, not only improves expected end-performance but also enhances the fairness of the rankings.

Moreover, in fairness intervals where the system's expected utility improves, it is possible that either
 fewer distracting items were included in the ranking passed to the generator or useful, previously
 overlooked items (which may have been considered random) were introduced due to the ranking
 perturbation. However, while higher utility paired with increased item-fairness (even within fairness
 intervals as low as [0.4, 0.6)) may seem advantageous, practitioners should exercise caution. This
 could result in compensating providers of items irrelevant to user requests, particularly in scenarios
 where content providers are rewarded for contributing to inference outcomes.

473 Machine-user browsing model. Developing more sophisticated machine-user browsing models 474 will result in more consistent and accurate evaluations of item-side fairness in RAG models, as the 475 exposure of each item is influenced by the attention allocated by the machine-user. Initial research 476 can draw inspiration from Liu et al. (2024), who found that LLMs tend to allocate more attention to 477 the beginning and end of a context, with less focus on the middle. This line of inquiry aligns with 478 the broader effort to create search engines tailored for machine-users (Salemi & Zamani, 2024b), 479 specifically focusing on fairer search engines in this context. It should involve studying how LMs 480 attend to each retrieved result within a context, analogous to how traditional search engine research 481 models human browsing behavior (Dupret & Piwowarski, 2008).

Measurement of string utility. In line with the recent call for evaluating various valid output strings (Zhang et al., 2024b), we recognize the need for a similar approach to better measure system utility across different rankings given. Recall that our experiments were designed to provide the generator with different rankings for the same query, leading to varied outputs. This approach is motivated by

the idea that items not appearing in the top positions of deterministic rankings may still hold value
and should be fairly considered by the system. In this context, the diverse outputs generated from
different rankings may still be valid. However, we currently rely on a single target output string for
comparison with predictions. Future work could focus on calculating the utility of diffuse predictions,
enabling a more nuanced evaluation.

Limitations. We acknowledge that the evaluation cost of fair RAG systems can be high due to repeated sampling and inference steps. However, in production, only a single ranking is sampled, minimizing the impact on system latency (Appendix E). Also, a limitation in our utility labeling is that it considers single items, while multiple items may yield contrasting utility gains. Despite this, the strong correlation between ranking quality and system effectiveness suggests this approach reasonably approximates item-worthiness for evaluating the impact of fair ranking on RAG systems.

497 **Conclusion**. This study highlights the impact of fair rankings on both the ranking and generation 498 quality of RAG systems. Through the extensive analysis, we show that fairer RAG systems not only 499 maintain high generation quality but can also outperform traditional RAG models, challenging the 500 notion of a strict tradeoff between fairness and effectiveness. Our findings provide valuable insights 501 for developing responsible and equitable RAG systems and pave the way for future research in fair 502 ranking and retrieval-augmented generation. In future work, we hope to extend this framework to 503 consider graded or missing judgments and exploring the different notions of fairness in RAG systems, ultimately advancing the field of trustworthy RAG systems research. 504

505 506

507

ETHICS STATEMENT

This study does not involve human subjects, harmful insights, or methodologies that raise ethical concerns. Additionally, there are no privacy, security, or legal issues associated with this work.
Instead, this research follows existing work on fair ranking (Ekstrand et al., 2022) that aims to promote equity while acknowledging the potential limitations of technical opertaionalizations of fairness.

513

514 REPRODUCIBILITY STATEMENT

515

516 The authors have made significant efforts to ensure the reproducibility of this research. 517 A well-documented anonymized code repository containing all necessary materials to 518 reproduce the experiments is available at https://anonymous.4open.science/r/ 519 fair-raq-iclr-anonymous-9C09. This repository describes a clear dataset source, a com-520 plete data generation pipeline using the LaMP benchmark (as detailed in §3.1 and Appendix B), and scripts for running the experiments. The code repository covers all stages of the research, including 521 deterministic retrieval computation, stochastic retrieval, RAG model implementation, Expected Expo-522 sure and Expected Utility evaluation, and normalization of metric values following the theoretical 523 proofs in Appendix C. For the part where the randomization is included a random seed is set to 42 524 across all experiments, and no sampling is used for language model decoding (as described in §4). 525 This approach was implemented not only to ensure reproducibility but also to limit the introduction 526 of stochasticity to the retrieval process, ensuring more reliable experiments. Implementations of nine 527 RAG models were based on the publicly available retrievers and generators which are referenced in 528 Appendix K. 529

530

531

References

- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-RAG: Learning to
 retrieve, generate, and critique through self-reflection. In *The Twelfth International Conference on Learning Representations*, 2024.
- Dara Bahri, Yi Tay, Che Zheng, Donald Metzler, and Andrew Tomkins. Choppy: Cut transformer for
 ranked list truncation. *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2020.

540 541 542 543	Aparna Balagopalan, Abigail Z. Jacobs, and Asia J. Biega. The role of relevance in fair ranking. In <i>Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , SIGIR '23, pp. 2650–2660. Association for Computing Machinery, 2023.
544 545 546	K. Balan, S. Agarwal, S. Jenni, A. Parsons, A. Gilbert, and J. Collomosse. Ekila: Synthetic media provenance and attribution for generative art. In 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), pp. 913–922. IEEE Computer Society, jun 2023.
547 548 549	Asia J Biega, Krishna P Gummadi, and Gerhard Weikum. Equity of attention: Amortizing individual fairness in rankings. In <i>The 41st international acm sigir conference on research & development in information retrieval</i> , pp. 405–414, 2018.
550 551 552 553	Sebastian Bruch, Shuguang Han, Michael Bendersky, and Marc Najork. A Stochastic Treatment of Learning to Rank Scoring Functions. In <i>Proceedings of the 13th International Conference on Web</i> <i>Search and Data Mining</i> , WSDM '20, pp. 61–69. Association for Computing Machinery, 2020.
554 555	Simon Caton and Christian Haas. Fairness in machine learning: A survey. <i>ACM Computing Surveys</i> , 2020.
556 557 558 559 560 561	Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language models, 2022.
562 563 564 565 566 567	Florin Cuconasu, Giovanni Trappolini, Federico Siciliano, Simone Filice, Cesare Campagnano, Yoelle Maarek, Nicola Tonellotto, and Fabrizio Silvestri. The power of noise: Redefining retrieval for rag systems. In <i>Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , SIGIR '24, pp. 719–729. Association for Computing Machinery, 2024.
568 569 570 571	Fernando Diaz, Bhaskar Mitra, Michael D. Ekstrand, Asia J. Biega, and Ben Carterette. Evaluating Stochastic Rankings with Expected Exposure. In <i>Proceedings of the 29th ACM International Conference on Information & Knowledge Management</i> , CIKM '20, pp. 275–284. Association for Computing Machinery, 2020.
572 573 574 575	Georges E. Dupret and Benjamin Piwowarski. A user browsing model to predict search engine click data from past observations. In <i>Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , SIGIR '08, pp. 331–338. Association for Computing Machinery, 2008. doi: 10.1145/1390334.1390392.
576 577 578	Michael D Ekstrand, Anubrata Das, Robin Burke, and Fernando Diaz. Fairness in information access systems. <i>Foundations and Trends</i> ® <i>in Information Retrieval</i> , 16(1-2):1–177, 2022.
579 580 581	Michael D. Ekstrand, Lex Beattie, Maria Soledad Pera, and Henriette Cramer. Not just algorithms: Strategically addressing consumer impacts in information retrieval. In <i>Advances in Information</i> <i>Retrieval</i> , pp. 314–335. Springer Nature Switzerland, 2024.
582 583 584 585 586	Shahul Es, Jithin James, Luis Espinosa Anke, and Steven Schockaert. RAGAs: Automated evaluation of retrieval augmented generation. In Nikolaos Aletras and Orphee De Clercq (eds.), <i>Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations</i> , pp. 150–158. Association for Computational Linguistics, 2024.
587 588 589 590	Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. Enabling large language models to generate text with citations. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pp. 6465–6488. Association for Computational Linguistics, 2023. doi: 10.18653/v1/2023.emnlp-main.398.
591 592 593	John Guiver and Edward Snelson. Bayesian inference for Plackett-Luce ranking models. In <i>Proceedings of the 26th Annual International Conference on Machine Learning</i> , pp. 377–384. ACM, 2009.

594 595 596	Emil Julius Gumbel. <i>Statistical theory of extreme values and some practical applications: a series of lectures</i> , volume 33. US Government Printing Office, 1954.
597 598 599	Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. Realm: Retrieval- augmented language model pre-training. In <i>Proceedings of the 37th International Conference on</i> <i>Machine Learning</i> , ICML'20. JMLR.org, 2020.
600 601	Peter Henderson, Xuechen Li, Dan Jurafsky, Tatsunori Hashimoto, Mark A. Lemley, and Percy Liang. Foundation models and fair use. <i>Journal of Machine Learning Research</i> , 24(400):1–79, 2023.
602 603 604 605	Sebastian Hofstätter, Jiecao Chen, Karthik Raman, and Hamed Zamani. Fid-light: Efficient and effective retrieval-augmented text generation. In <i>Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pp. 1437–1447, 2023.
606 607 608	Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. Unsupervised dense information retrieval with contrastive learning. <i>Transactions on Machine Learning Research</i> , 2022. ISSN 2835-8856.
609 610 611 612	Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. Atlas: Few-shot learning with retrieval augmented language models. <i>Journal of Machine Learning Research</i> , 24(251):1–43, 2023.
613 614 615 616 617	Thomas Jaenich, Graham McDonald, and Iadh Ounis. Fairness-aware exposure allocation via adaptive reranking. In <i>Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , SIGIR '24, pp. 1504–1513. Association for Computing Machinery, 2024.
618 619 620 621	Meng Jiang, Keqin Bao, Jizhi Zhang, Wenjie Wang, Zhengyi Yang, Fuli Feng, and Xiangnan He. Item-side fairness of large language model-based recommendation system. In <i>Proceedings of the ACM on Web Conference 2024</i> , WWW '24, pp. 4717–4726. Association for Computing Machinery, 2024.
622 623	Diane Kelly, Jaime Arguello, and Robert Capra. Nsf workshop on task-based information search systems. <i>SIGIR Forum</i> , 47(2):116–127, 2013. ISSN 0163-5840.
625 626 627	Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. Generalization through memorization: Nearest neighbor language models. In <i>International Conference on Learning Representations</i> , 2020.
628 629	To Eun Kim, Alireza Salemi, Andrew Drozdov, Fernando Diaz, and Hamed Zamani. Retrieval- enhanced machine learning: Synthesis and opportunities. <i>arXiv preprint arXiv:2407.12982</i> , 2024.
630 631 632 633	Ishita Kumar, Snigdha Viswanathan, Sushrita Yerra, Alireza Salemi, Ryan A. Rossi, Franck Dernon- court, Hanieh Deilamsalehy, Xiang Chen, Ruiyi Zhang, Shubham Agarwal, Nedim Lipka, and Hamed Zamani. Longlamp: A benchmark for personalized long-form text generation, 2024.
634 635 636 637 638 639	Patrick S. H. Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive NLP tasks. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), <i>Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual</i> , 2020.
640 641 642	Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts. <i>Transactions of the Association for Computational Linguistics</i> , 12:157–173, 2024.
643 644 645	Lingjuan Lyu, C Chen, and J Fu. A pathway towards responsible ai generated content. In <i>IJCAI</i> , pp. 7033–7038, 2023.
646 647	Chris J. Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. In <i>International Conference on Learning Representations</i> , 2017.

648 Rishabh Mehrotra, Ashton Anderson, Fernando Diaz, Amit Sharma, Hanna Wallach, and Emine 649 Yilmaz. Auditing search engines for differential satisfaction across demographics. In Proceedings 650 of the 26th International Conference on World Wide Web Companion, WWW '17 Companion, pp. 651 626-633. International World Wide Web Conferences Steering Committee, 2017. 652 Alistair Moffat and Justin Zobel. Rank-biased precision for measurement of retrieval effectiveness. 653 ACM Transactions on Information Systems (TOIS), 27(1):1–27, 2008. 654 655 Abhiman Neelakanteswara, Shreyas Chaudhari, and Hamed Zamani. Rags to style: Personalizing 656 llms with style embeddings. In Proceedings of the 1st Workshop on Personalization of Generative 657 AI Systems (PERSONALIZE 2024), pp. 119–123, 2024. 658 Harrie Oosterhuis. Computationally efficient optimization of plackett-luce ranking models for 659 relevance and fairness. In Proceedings of the 44th International ACM SIGIR Conference on 660 Research and Development in Information Retrieval, pp. 1023–1032, 2021. 661 662 Harrie Oosterhuis. Learning-to-rank at the speed of sampling: Plackett-luce gradient estimation with minimal computational complexity. In Proceedings of the 45th International ACM SIGIR 663 Conference on Research and Development in Information Retrieval, pp. 2266–2271, 2022. 664 665 Robin L Plackett. The analysis of permutations. Journal of the Royal Statistical Society Series C: 666 Applied Statistics, 24(2):193–202, 1975. 667 Amifa Raj and Michael D Ekstrand. Comparing fair ranking metrics. arXiv preprint 668 arXiv:2009.01311, 2020. 669 670 Stephen Robertson, S. Walker, S. Jones, M. M. Hancock-Beaulieu, and M. Gatford. Okapi at trec-3. 671 In Proceedings of the Third Text REtrieval Conference, TREC-3, pp. 109–126. Gaithersburg, MD: 672 NIST, 1995. 673 Dongyu Ru, Lin Qiu, Xiangkun Hu, Tianhang Zhang, Peng Shi, Shuaichen Chang, Jiayang Cheng, 674 Cunxiang Wang, Shichao Sun, Huanyu Li, et al. Ragchecker: A fine-grained framework for 675 diagnosing retrieval-augmented generation. arXiv preprint arXiv:2408.08067, 2024. 676 677 Jon Saad-Falcon, Omar Khattab, Christopher Potts, and Matei Zaharia. ARES: An automated 678 evaluation framework for retrieval-augmented generation systems. In Kevin Duh, Helena Gomez, 679 and Steven Bethard (eds.), Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long 680 Papers), pp. 338–354. Association for Computational Linguistics, 2024. 681 682 Alireza Salemi and Hamed Zamani. Evaluating retrieval quality in retrieval-augmented generation. 683 In Proceedings of the 47th International ACM SIGIR Conference on Research and Development in 684 Information Retrieval, pp. 2395–2400, 2024a. 685 Alireza Salemi and Hamed Zamani. Towards a search engine for machines: Unified ranking for 686 multiple retrieval-augmented large language models. In Proceedings of the 47th International 687 ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '24, pp. 688 741–751. Association for Computing Machinery, 2024b. doi: 10.1145/3626772.3657733. 689 690 Alireza Salemi, Surya Kallumadi, and Hamed Zamani. Optimization methods for personalizing 691 large language models through retrieval augmentation. In Proceedings of the 47th International 692 ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '24, pp. 693 752–762. Association for Computing Machinery, 2024a. doi: 10.1145/3626772.3657783. 694 Alireza Salemi, Sheshera Mysore, Michael Bendersky, and Hamed Zamani. LaMP: When large 695 language models meet personalization. In Proceedings of the 62nd Annual Meeting of the Associ-696 ation for Computational Linguistics (Volume 1: Long Papers), pp. 7370–7392. Association for 697 Computational Linguistics, 2024b. 698 Piotr Sapiezynski, Wesley Zeng, Ronald E Robertson, Alan Mislove, and Christo Wilson. Quantifying 699 the impact of user attention fair group representation in ranked lists. In Companion Proceedings 700 of The 2019 World Wide Web Conference, WWW '19, pp. 553–562. Association for Computing 701 Machinery, 2019.

702 703 704	Tefko Saracevic. The Notion of Relevance in Information Science: Everybody knows what relevance is. But, what is it really? Morgan & Claypool Publishers, 2016.
705 706	Robik Shrestha, Yang Zou, Qiuyu Chen, Zhiheng Li, Yusheng Xie, and Siqi Deng. Fairrag: Fair human generation via fair retrieval augmentation. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pp. 11996–12005, 2024.
708	Ashudeen Singh and Thorsten Joachims, Fairness of exposure in rankings. In Proceedings of the 24th
709	ACM SIGKDD international conference on knowledge discovery & data mining, pp. 2219–2228,
710	2018.
711	All the O's the TTL store Testing Dation for Conference is setting Te Date to Conference is setting to Date the
712 713	Ashudeep Singh and Thorsten Joachims. Policy learning for fairness in ranking. In Proceedings of the 33rd International Conference on Neural Information Processing Systems, number 487, pp.
714	5420–5450. Curran Associates Inc., 2019.
715 716	Ke Yang and Julia Stoyanovich. Measuring fairness in ranked outputs. In <i>Proceedings of the 29th international conference on scientific and statistical database management</i> , pp. 1–6, 2017.
717	
718	Hamed Zamani and Michael Bendersky. Stochastic rag: End-to-end retrieval-augmented generation
719	through expected utility maximization. In <i>Proceedings of the 47th International ACM SIGIR</i>
720	Conference on Research and Development in Information Retrieval, SIGIR 724, pp. 2641–2646.
721	Association for Computing Machinery, 2024. doi: 10.1145/3626772.3657923.
722	Hamed Zamani, Fernando Diaz, Mostafa Dehghani, Donald Metzler, and Michael Bendersky,
723	Retrieval-enhanced machine learning. In <i>Proceedings of the 45th Annual International ACM</i>
724	SIGIR Conference on Research and Development in Information Retrieval, 2022.
725	Meiler Zehliler Franzesse Denski Casles Castille Sam Heilen Mahamad Masshed and Disarda
726	Baeza Vates Ea* ir: A fair top k ranking algorithm In <i>Proceedings of the 2017 ACM on</i>
727	Conference on Information and Knowledge Management, pp. 1569–1578, 2017.
728	Huimin Zeng, Zhenrui Yue, Qian Jiang, and Dong Wang. Federated recommendation via hybrid
729	retrieval augmented generation, 2024.
731	
732	Hengran Zhang, Ruqing Zhang, Jiafeng Guo, Maarten de Rijke, Yixing Fan, and Xueqi Cheng.
733	ACM SIGIR Conference on Research and Development in Information Retrieval SIGIR '24 pp
734	1941–1951. Association for Computing Machinery, 2024a.
735	Virging Zhang, Avi Schwarzschild Nicholas Carlini, Zies Kalter, and Danhus Involite. Foreing
736 737	diffuse distributions out of language models. <i>arXiv preprint arXiv:2404.10859</i> , 2024b.
738	
739	
740	
741	
742	
743	
744	
745	
746	
747	
748	
749	
750	
751	
752	
753	
754	

A NOTATION

758		
759	Notation	Description
760	\overline{r}	input instance
761	$\frac{\omega}{v}$	output target
762	$\phi_a(x)$	query generation function
763	q	query returned by $\phi_a(x)$
764	\dot{d}	retrieval item (document)
765	С	stored retrievable items (corpus)
766	n	the number of d's in C
67	m	the number of useful items in C
68		
769	$\mathcal{R}(q, \mathbf{C})$	deterministic retriever
70	L	ranked list returned by $\mathcal{R}(q, C)$
771	$\mathbf{s} \in \mathbb{K}^{n}$	retrieval scores returned by $\mathcal{K}(q, \mathbb{C})$
772	IN k	sampling size for stochastic sampling the number of d 's to retrieve per ranking
773	N S(s N L)	stochastic ranking sampler
774	σ	set of N sampled rankings returned by $S(s \ N \ k)$
775	π	sampled ranking $\in \sigma$
76	$\overset{n}{\phi}_{n}(x,\pi)$	prompt generation function
77	$\frac{\overline{x}}{\overline{x}}$	prompt returned by $\phi_n(x,\pi)$
78	$\mathcal{G}(\overline{x})$	language model
70	\hat{y}	predicted output from $\mathcal{G}(\overline{x})$
790		
701	wor(d x)	worthiness of an item d given an input x
782	$\mu_f(\sigma)$	fairness metric of rankings
783	$\mu_r(\sigma)$	relevance metric of rankings
79.4	$\mu_u(y,y)$	string utility metric
795	$\epsilon \in \mathbb{R}^n$	target exposure of all items in C
700	e e m	target exposure of an items in C
700		Table 2: Notation
700		
700		
790		
701		
700		
32		
30		
94		
190		
96		
/9/		
798		
799		
300		
301		
802		

810 B LABELING PROCEDURE 811

1: D	$= \{(x_1, y_1), (x_2, y_2), \cdots, (x_n)\}$	(T, y_T) > dataset of size T
2: fo	$(x_i, y_i) \in \mathcal{D}$ do	
3:	$u_i \leftarrow \mu_u(y_i, \mathcal{G}(x_i))$	▷ string utility of a baseline model without augmentation
4:	for $d_j \in \mathbf{C}$ do	
5:	$\hat{y}_j \leftarrow \mathcal{G}(\phi_p(x_i, d_j))$	
6:	$u_j \leftarrow \mu_u(y_i, \hat{y}_j)$	\triangleright string utility of a generator augmented with one item
7:	$\delta_j \leftarrow (u_j - u_i)$	\triangleright marginal utility gained from the augmentation
8:	$wor(a_j x_i) \leftarrow 0$ if $\delta > 0$ then	b hinery decision of item worthings by the utility gain
9:	$uor(d x_i) \leftarrow 1$	binary decision of item-worthiness by the utility-gain
10.	end if	
12:	end for	
13: en	d for	

С NORMALIZATION OF METRICS

C.1 NORMALIZATION OF EE-D

The disparity measure EE-D should be normalized by its true upper and lower bound.

Theorem 1. $\|\epsilon\|_2^2 \in [0, \|\bar{\epsilon}\|_2^2]$, where $\bar{\epsilon}$ is an exposure vector derived from any deterministic ranking.

Proof. The lower bound is achieved by a uniform random policy. Each item d will have exposure of $\frac{1}{2}$. However, it is reasonable to assume that it is approximately 0, since the size of most of the retrieval corpus is very large. Also, it is common that the corpus consists of majority of non-relevant items. The implication is that, for the optimal exposure, since the n-m non-relevant items are shuffled amongst themselves, each will have an expected exposure of close to 0. Thus, assuming large n and relatively small m,

$$\frac{1}{n} < \frac{1}{n-m} \approx 0 \tag{8}$$

For upper bound, recall that the ϵ is computed based on samples of rankings from a stochastic policy. For relevance-independent browsing models (e.g., MU), all rankings $\pi \in S_n$ have identical exposure norms $\|\epsilon^{\pi}\|_{2}^{2}$, where ϵ^{π} is the exposure of items for ranking π sampled from the stochastic policy. Then,

$$\|\epsilon\|_{2}^{2} = \|\mathbb{E}_{\pi}[\epsilon^{\pi}]\|_{2}^{2} \tag{9}$$

$$\leq \mathbb{E}_{\pi}[\|\epsilon^{\pi}\|_{2}^{2}] \tag{10}$$

$$= \mathbb{E}_{\pi}[\|\bar{\epsilon}\|_{2}^{2}] = \|\bar{\epsilon}\|_{2}^{2}$$
(11)

Corollary 1. For machine user browsing model with top-k consumption and equal attention to the *top items*, $\|\epsilon\|_{2}^{2} \in [0, k]$

With MU, the upper bound of EE-D, $\|\bar{\epsilon}\|_2^2$ becomes $\sum_{i=1}^n MU(i)^2 = k$. Therefore, per query q, we calculate a normalized EE-D Ē

$$\overline{\text{E-D}}_q = \|\epsilon\|_2^2 / k \quad \in [0,1] \tag{12}$$

C.2 NORMALIZATION OF EE-R

The ranking quality measure EE-R should be normalized by its true upper and lower bound.

Theorem 2. $\langle \epsilon, \epsilon^* \rangle \in [0, \|\epsilon^*\|_2^2]$

> *Proof.* The lower bound is achieved when ϵ becomes ϵ^- , which is an exposure vector of the worst case ranking π^- (permutations that rank all non-relevant items above relevant items). Given the assumption made from equation 8 and, C^+ and C^- , which are set of indices of relevant and nonrelevant items, respectively,

$$\langle \epsilon^-, \epsilon^* \rangle = \sum_{i=1}^n \epsilon_i^- \epsilon_i^* \tag{13}$$

$$\approx \sum_{i \in \mathbb{C}^+} 0\epsilon_i^* + \sum_{i \in \mathbb{C}^-} \epsilon_i^- 0 = 0$$
 (from 8) (14)

Intuitively, the upper bound is achieved when ϵ becomes ϵ^* , thus $\langle \epsilon^*, \epsilon^* \rangle = \|\epsilon^*\|_2^2$. Alternatively, we can show that any convex combination of optimal rankings will have a $\langle \epsilon, \epsilon^* \rangle = \|\epsilon^*\|_2^2$.

P18 Let ϵ_d^* be the exposure of a relevant items, S_n^* be the set of all optimal rankings, $w_{\pi'}$ be the weight on $\pi \in S_n^*$ such that $\sum_{\pi \in S_n^*} w_{\pi'} = 1$, and ϵ' be the exposure vector associated with π' .

921
922
922
922

$$\langle \epsilon, \epsilon^* \rangle = \sum_{i=1}^n \epsilon_i \epsilon_i^* = \sum_{i \in C^+} \epsilon_i \epsilon_i^*$$
 (from 8) (15)

$$\stackrel{i=1}{= \epsilon_{d}^{*} \sum \epsilon_{i}} \epsilon_{i} \qquad (equal exposure principle) \qquad (16)$$

$$=\epsilon_d^* \sum_{\pi' \in S^*} w_{\pi'} \sum_{i \in C^+} \epsilon_i'$$
(18)

$$\leq \epsilon_d^* \sum_{\pi' \in S_n^*} w_{\pi'}(m\epsilon_d^*) \qquad (\text{since } \epsilon' \text{ is optimal}) \tag{19}$$

$$= \epsilon_d^*(m\epsilon_d^*) \qquad (\text{since } \sum_{\pi \in S_n^*} w_{\pi'} = 1) \qquad (20)$$

$$= \sum_{i \in C^+} \epsilon_i^* \epsilon_i^* = \|\epsilon^*\|_2^2$$
(21)

Corollary 2. For machine user browsing model with top-k consumption and equal attention to the top items, the bound depends on m and k. If $m \le k$, $\langle \epsilon, \epsilon^* \rangle \in [0, m + \frac{(k-m)^2}{n-m}]$. If m > k, $\langle \epsilon, \epsilon^* \rangle \in [0, \frac{k^2}{m}]$

943 With MU, the upper bound of EE-R can be calculated by equation 6. 944 If $m \le k$, $\|\epsilon^*\|_2^2$ becomes

$$\sum_{i=C^{+}} 1^{2} + \sum_{i=C^{-}} \left(\frac{k-m}{n-m}\right)^{2} = m + \frac{(k-m)^{2}}{n-m}$$
(22)

If m > k, $\|\epsilon^*\|_2^2$ becomes

$$\sum_{i=C^{+}} \left(\frac{k}{m}\right)^{2} + \sum_{i=C^{-}} 0^{2} = \frac{k^{2}}{m}$$
(23)

Therefore, depending on m and k, per query q, we calculate a normalized EE-R

$$\overline{\text{EE-R}}_q = \begin{cases} \langle \epsilon, \epsilon^* \rangle / (m + \frac{(k-m)^2}{n-m}) & (m \le k) \\ \langle \epsilon, \epsilon^* \rangle / (\frac{k^2}{m}) & (m > k) \end{cases} \in [0,1]$$
(24)

C.3 NORMALIZATION OF EU

Theoretically obtaining a true upper bound of a utility of a RAG model is challenging. Therefore, we approximate the true upper bound by the maximum of the empirically obtained utilities given a fixed RAG model $\langle S, G \rangle$ with a stochastic ranking sampler.

963 964 Recall that the string utility $u_{\pi} = \mu_u(y, \mathcal{G}(\phi_p(x, \pi)))$ is a utility of a RAG model with a sampled 965 ranking $\pi \in \sigma$. Let σ_{α} denote a sample of rankings with temperature parameter set to α . Also, let σ^* 966 denote the set of sampled permutations (rankings) from the oracle stochastic retriever, a policy that 967 always places relevant items above non-relevant items; thus the oracle generates m!(n-m)! number 967 of unique optimal permutations.

To obtain an approximated upper bound of the utility u_{max} , when the runs of the experiments were run with $\alpha = (1, 2, 4, 8)$, we take the maximum over all samples,

$$u_{max} = max \left(\{u_{\pi}\}_{\pi \in \sigma_1} \cup \{u_{\pi}\}_{\pi \in \sigma_2} \cup \{u_{\pi}\}_{\pi \in \sigma_4} \cup \{u_{\pi}\}_{\pi \in \sigma_8} \cup \{u_{\pi}\}_{\pi \in \sigma^*} \right)$$
(25)

951 952 953

959

971

949 950

924

934 935 936

With u_{max} , we max-normalize the EU. Since $\frac{1}{N} \sum_{\pi \in \sigma} \frac{u_{\pi}}{u_{max}}$ is the same as $(\frac{1}{N} \sum_{\pi \in \sigma} u_{\pi})/u_{max}$, per query q, we get a normalized EU

$$\overline{\mathrm{EU}}_q = \frac{\mathrm{EU}_q}{u_{max}} \quad \in [0, 1] \tag{26}$$

Normalization of EU is done to get the percentage of closeness to the optimal utility as all the utility
values are scaled relative to the maximum value. In other words, the normalized EU value indicates
how close the EU is to the maximum utility that the RAG system can get to.

This is straightforward for *higher-the-better* metrics, such as ROUGE and Accuracy. However, for
 lower-the-better metrics such as MAE, we convert the scores to *higher-the-better* by subtracting the
 scores from the true metric upper bound. This allows us to perform the same normalization operation
 and have the same interpretation of the normalized metric.



(c) generation quality of RAG. The left column is the results on the LaMP task 1 and the right column is on the LaMP task 4, each corresponding to a classification and text generation task, respectively. RAG model is configured with Contriever and Flan-T5-Base for all six figures. Each data point 1062 represents the value of the associated metric for one query. 1063

1064

1026

D

When the three types of plots are observed together, we can infer some interesting observations. In 1065 general, we see positive relationships between increasing α and both ranking quality and utility. This 1066 implies that we can generally expect a tradeoff between both fairness and ranking quality, as well as 1067 between fairness and utility. 1068

However, we can expect some edge cases. For instance, in LaMP-1 (left column of the figure), 1069 difference in EU when $\alpha = 4$ and $\alpha = 8$ is not large, and we can see that even the lower quartile of 1070 the utility is increased when α is set to 4 (Figure 6c). This can possibly mean that in LaMP-1, the 1071 RAG model can maintain considerable utility when the disparity is roughly in the range of [0.6, 0.8]. 1072 Similar observation can be made for the LaMP-4 (right column of the figure), except that the $\overline{\text{EE-R}}$ is 1073 higher when $\alpha = 4$ than when $\alpha = 8$ (Figure 6b). This indicates that the deterministic retriever does 1074 not always provide the maximum ranking quality, and the retriever can sometimes provide higher 1075 ranking quality by being more fair, ultimately maintaining considerable or higher utility (similar EU 1076 when $\alpha=4$ and when $\alpha=8$) at the same time. 1077

With this preliminary observations in mind, we could delve deeper into the relationships between 1078 fairness, ranking quality, and utility, by visualizing and quantifying the combined results from all the 1079 runs ($\alpha = 1, 2, 4, 8$).

E INTEGRATING FAIR RANKINGS INTO A RAG SYSTEM IN A PRODUCTION ENVIRONMENT



Figure 7: Designing a RAG system that incorporates stochastic fair rankings involves using stochastic sampling, where N can be set to 1 to provide a single ranking to the language model. This can result in a different ranking compared to the one in Figure 1, exposing d_1 and d_3 to the language model. This paper is focusing on the evaluation of this system as depicted in the Figure 2. The query and prompt generators are omitted in the figure for brevity.

¹¹³⁴ F FAIRNESS VS. RANKING QUALITY

1136



F.1 VISUALIZATION OF EE-D VS. EE-R OF FLANT5SMALL







1242 F.3 VISUALIZATION OF EE-D VS. EE-R OF FLANT5XXL

1296 F.4 QUANTIFICATION OF FAIRNESS-RANKING QUALITY TRADEOFF

1208										
1230		slope↓ / AUC↑								
1200	Task		FlanT5Small			FlanT5Base			FlanT5XXL	
1233	Iask	BM25	SPLADE	Contriever	BM25	SPLADE	Contriever	BM25	SPLADE	Contriever
1300	LaMP-1	0.2113/0.2546	0.2396 / 0.3147	0.2358 / 0.2102	0.1546/0.3574	0.2252 / 0.4043	0.1834 / 0.3594	0.1358 / 0.3390	0.2409 / 0.4130	0.1945 / 0.4009
1000	LaMP-2	0.1599 / 0.3079	0.1863 / 0.3330	0.2072 / 0.3372	0.1665 / 0.3699	0.1899 / 0.3918	0.2269 / 0.4098	0.1651/0.4161	0.1834 / 0.4222	0.2029 / 0.4328
1301	LaMP-3	0.0309 / 0.3820	0.0353 / 0.3970	0.0501/0.3911	0.0403 / 0.4731	0.0290 / 0.4760	0.0555 / 0.4745	0.0798 / 0.3169	0.1010 / 0.3483	0.1338 / 0.3502
1001	LaMP-4	0.1271/0.3425	0.1702 / 0.3741	0.1858 / 0.3793	0.1054 / 0.3812	0.1715 / 0.4236	0.1837 / 0.4303	0.0882 / 0.3669	0.1286 / 0.3973	0.1363 / 0.4035
1302	LaMP-5	0.1058 / 0.3396	0.1014 / 0.3386	0.1031/0.3356	0.0946 / 0.3477	0.1144 / 0.3597	0.1157 / 0.3533	0.2072 / 0.4026	0.2071 / 0.4100	0.2086 / 0.4040
1302	LaMP-6	0.1194 / 0.4507	0.1036 / 0.4477	0.1090 / 0.4491	0.1456 / 0.4880	0.1264 / 0.4967	0.1438 / 0.4942	0.1615 / 0.5685	0.1554 / 0.5813	0.1467 / 0.5747
1303	LaMP-7	0.0962 / 0.4984	0.1178 / 0.4976	0.0979 / 0.4894	0.0744 / 0.3811	0.1173 / 0.3908	0.0924 / 0.3926	0.1174 / 0.4800	0.0216 / 0.4784	0.0627 / 0.4885

Table 3: Values on the left are the gradient of a linear line fit to the data points where x-axis is $\overline{\text{EE-D}}$ and y-axis is $\overline{\text{EE-R}}$. Higher the value, stronger the tradeoff between fairness and ranking quality. Values on the right are the DR-AUC on the disparity-ranking quality ($\overline{\text{EE-D}}$ Vs. $\overline{\text{EE-R}}$) curve. Higher the value, stronger the ranking quality, given consistent tradeoff between fairness and relevance.

¹³⁵⁰ G RANKING QUALITY VS. UTILITY OF RAG MODELS



G.2 QUANTIFICATION OF THE RELATIONSHIP BETWEEN RANKING QUALITY AND UTILITY

	slope⊥ / AUC↑	slope⊥ / AUC↑	slope⊥ / AUC↑	slope⊥ / AUC↑	slope⊥ / AUC↑	slope⊥ / AUC↑	slope⊥ / AUC↑	slope. / AUC↑	slope⊥ / AUC↑
T1.	1.1.	BM25	1 1		SPLADE			Contriever	1.1.
Task	FlanT5Small	FlanT5Base	FlanT5XXL	FlanT5Small	FlanT5Base	FlanT5XXL	FlanT5Small	FlanT5Base	FlanT5XXL
LaMP-1	0.3813 / 0.3705	0.8250 / 0.7491	0.7612/0.7044	0.1434 / 0.3760	0.7783 / 0.7382	0.7454 / 0.7330	0.1864 / 0.3514	0.7661/0.7304	0.7089/0.7166
LaMP-2	0.5543 / 0.4130	0.2858 / 0.3218	0.2646 / 0.4495	0.5539/0.3941	0.3310/0.3407	0.1979 / 0.4787	0.6908 / 0.4567	0.2325 / 0.3600	0.3342 / 0.4772
LaMP-3	0.2500 / 0.9153	0.2134 / 0.9204	0.1633 / 0.8974	0.2544 / 0.9124	0.2061 / 0.9203	0.1737 / 0.9009	0.2538 / 0.9125	0.2025 / 0.9205	0.1568 / 0.8977
LaMP-4	0.2708 / 0.2711	0.2881 / 0.2679	0.1947 / 0.3720	0.2638/0.2817	0.3272 / 0.2942	0.2193 / 0.3961	0.2264 / 0.2751	0.2986 / 0.2887	0.2190/0.3926
LaMP-5	-0.0125 / 0.3723	-0.0307 / 0.4888	0.1673/0.5591	-0.0035 / 0.3849	-0.0160/0.4773	0.2044 / 0.5570	0.0090 / 0.3823	-0.0474 / 0.4746	0.1590/0.5387
LaMP-6	0.2575 / 0.3883	0.2690 / 0.3953	0.3078 / 0.5445	0.2522 / 0.3737	0.2735 / 0.3960	0.2780 / 0.5372	0.2790 / 0.3783	0.2620 / 0.3891	0.3053 / 0.5356
LaMP-7	0.3229 / 0.5240	0.0889 / 0.6391	0.0157 / 0.6082	0.2781 / 0.5177	0.2029 / 0.6435	-0.0623 / 0.5601	0.2713/0.5101	0.0788 / 0.6266	-0.0466 / 0.5723

<sup>Table 4: Values on the left are the gradient of a linear line fit to the data points where x-axis is EE-R
and y-axis is EU. Higher the value, stronger the tradeoff between retrieval quality and generation
quality. Values on the right are the RU-AUC on the ranking quality-utility (EE-RVs. EU) curve.
Higher the value, stronger the general end-performance of a RAG model when every level of relevance is considered.</sup>

1404 H ITEM-FAIRNESS VS. UTILITY OF RAG MODELS

H.1 VISUALIZATION OF EE-D VS. EU



1438 H.2 QUANTIFICATION OF FAIRNESS-UTILITY TRADEOFF

	slope↓ / AUC↑	slope↓ / AUC↑	slope↓ / AUC↑	slope↓ / AUC↑	slope↓ / AUC↑	slope↓ / AUC↑	slope↓ / AUC↑	slope↓ / AUC↑	slope↓ / AUC↑
Tack		BM25			SPLADE			Contriever	
Task	FlanT5Small	FlanT5Base	FlanT5XXL	FlanT5Small	FlanT5Base	FlanT5XXL	FlanT5Small	FlanT5Base	FlanT5XXL
LaMP-1	0.0693 / 0.1994	0.2254/0.6110	0.1673 / 0.5688	0.0851/0.2519	0.1998 / 0.6644	0.2385 / 0.6362	0.0456/0.1866	0.2016 / 0.6248	0.2413/0.6061
LaMP-2	0.1295 / 0.2561	0.0740 / 0.2717	0.0600 / 0.4294	0.1119/0.2722	0.0637 / 0.2786	0.0626 / 0.4582	0.1683 / 0.3082	0.0870/0.3139	0.1002 / 0.4338
LaMP-3	0.0259 / 0.8738	-0.0012 / 0.8964	0.0101 / 0.8706	0.0118 / 0.8744	0.0014 / 0.9003	0.0280 / 0.8786	0.0198 / 0.8685	0.0063 / 0.9022	0.0264 / 0.8765
LaMP-4	0.0606 / 0.2178	0.0789 / 0.2282	0.0734 / 0.3622	0.0937 / 0.2408	0.1245 / 0.2594	0.0895 / 0.3802	0.1016 / 0.2429	0.1209 / 0.2595	0.1051/0.3784
LaMP-5	0.0293 / 0.3995	0.0445 / 0.4857	0.0363 / 0.5453	0.0437 / 0.3989	0.0266 / 0.4763	0.0320 / 0.5466	0.0533 / 0.3930	0.0300 / 0.4943	0.0432 / 0.5495
LaMP-6	0.1177 / 0.3698	0.0753 / 0.3861	0.0521/0.5535	0.0806 / 0.3629	0.0623 / 0.3888	0.0561 / 0.5558	0.0823 / 0.3580	0.0597 / 0.3849	0.0616 / 0.5468
LaMP-7	0.0215 / 0.4928	0.0392 / 0.6212	0.0854 / 0.6068	0.0718/0.4901	0.0256 / 0.6094	-0.0668 / 0.5714	0.0324 / 0.4839	-0.0151 / 0.6101	-0.0114 / 0.5809

<sup>Table 5: Values on the left are the gradient of a linear line fit to the data points where x-axis is EE-D
and y-axis is EU. Higher the value, stronger the tradeoff between item-fairness and generation quality.
Values on the right are the DU-AUC on the disparity-utility (EE-D Vs. EU) curve. Higher the value, stronger the general end-performance of a RAG model when every level of fairness is considered.</sup>

1458 I PERFORMANCE OF RAG MODELS WITH VARYING FAIRNESS LEVELS

	Fairness Intervals							
Model (baseline utility)	[0.0, 0.2)	[0.2, 0.4)	[0.4, 0.6)	[0.6, 0.8)	[0.8, 1.0)			
BM25+FlanT5Small (0.308)	-0.12	-0.13	-0.18	-0.02	-0.15			
BM25+FlanT5Base (0.670)	-0.20	-0.04	-0.08	-0.05	-0.02			
BM25+FlanT5XXL (0.531)	-0.07	+0.03	+0.02	+0.06	+0.11			
SPLADE+FlanT5Small (0.241)	-0.03	-0.22	+0.19	-0.04	+0.14			
SPLADE+FlanT5Base (0.646)	-0.15	+0.06	+0.08	0.00	+0.03			
SPLADE+FlanT5XXL (0.671)	-0.18	-0.16	+0.05	+0.02	+0.01			
Contriever+FlanT5Small (0.286)	-0.08	-0.29	-0.06	+0.03	-0.14			
Contriever+FlanT5Base (0.637)	-0.16	+0.05	-0.06	+0.03	0.00			
Contriever+FlanT5XXL (0.651)	-0.19	-0.04	-0.11	+0.03	0.00			

Table 6: LaMP-1

	Fairness Intervals							
Model (baseline utility)	[0.0, 0.2)	[0.2, 0.4)	[0.4, 0.6)	[0.6, 0.8)	[0.8, 1.0)			
BM25+FlanT5Small (0.274)	-0.13	+0.02	-0.04	-0.04	+0.03			
BM25+FlanT5Base (0.223)	-0.01	+0.05	+0.04	+0.03	+0.13			
BM25+FlanT5XXL (0.310)	+0.05	+0.02	+0.24	+0.13	+0.18			
SPLADE+FlanT5Small (0.209)	-0.06	+0.10	+0.06	+0.05	+0.14			
SPLADE+FlanT5Base (0.238)	-0.02	+0.04	+0.04	+0.05	+0.09			
SPLADE+FlanT5XXL (0.472)	-0.05	-0.14	+0.12	0.00	-0.02			
Contriever+FlanT5Small (0.318)	-0.15	+0.05	-0.05	-0.04	+0.06			
Contriever+FlanT5Base (0.302)	-0.07	+0.02	+0.01	+0.02	+0.05			
Contriever+FlanT5XXL (0.356)	0.00	0.00	+0.12	+0.12	+0.15			

Table 7: LaMP-2

	Fairness Intervals						
Model (baseline utility)	[0.0, 0.2)	[0.2, 0.4)	[0.4, 0.6)	[0.6, 0.8)	[0.8, 1.0)		
BM25+FlanT5Small (0.886)	-0.03	+0.01	-0.01	-0.03	-0.01		
BM25+FlanT5Base (0.907)	-0.02	+0.01	0.00	-0.01	-0.04		
BM25+FlanT5XXL (0.859)	-0.02	+0.07	+0.01	+0.01	-0.02		
SPLADE+FlanT5Small (0.847)	+0.01	+0.05	+0.03	+0.01	+0.04		
SPLADE+FlanT5Base (0.902)	-0.01	+0.03	0.00	-0.01	-0.02		
SPLADE+FlanT5XXL (0.864)	-0.02	+0.08	0.00	+0.01	0.00		
Contriever+FlanT5Small (0.876)	-0.02	0.00	-0.02	0.00	-0.01		
Contriever+FlanT5Base (0.894)	0.00	+0.03	0.00	+0.02	0.00		
Contriever+FlanT5XXL (0.865)	-0.02	+0.07	+0.01	-0.02	+0.01		

Table 8: LaMP-3

		Fa	irness Interv	als	
Model (baseline utility)	[0.0, 0.2)	[0.2, 0.4)	[0.4, 0.6)	[0.6, 0.8)	[0.8, 1.0)
BM25+FlanT5Small (0.217)	-0.06	0.00	+0.02	+0.01	0.00
BM25+FlanT5Base (0.223)	-0.06	0.00	+0.03	+0.01	+0.02
BM25+FlanT5XXL (0.322)	-0.05	+0.11	+0.03	+0.03	+0.05
SPLADE+FlanT5Small (0.235)	-0.07	-0.01	+0.02	+0.03	+0.02
SPLADE+FlanT5Base (0.268)	-0.10	-0.03	+0.02	0.00	+0.02
SPLADE+FlanT5XXL (0.342)	-0.06	+0.09	+0.05	+0.03	+0.04
Contriever+FlanT5Small (0.254)	-0.09	-0.02	0.00	+0.01	0.00
Contriever+FlanT5Base (0.268)	-0.10	-0.02	+0.01	0.00	+0.01
Contriever+FlanT5XXL (0.367)	-0.09	+0.06	+0.01	+0.01	+0.03

Table 9: LaMP-4

		Fo	rnace Interv	o1c	
		Fa	inness miterv	a15	
Model (baseline utility)	[0.0, 0.2)	[0.2, 0.4)	[0.4, 0.6)	[0.6, 0.8)	[0.8, 1.0)
BM25+FlanT5Small (0.343)	+0.01	+0.12	+0.06	+0.04	+0.06
BM25+FlanT5Base (0.507)	-0.04	-0.04	-0.01	-0.02	-0.01
BM25+FlanT5XXL (0.508)	-0.03	+0.16	+0.02	+0.02	0.00
SPLADE+FlanT5Small (0.378)	-0.03	+0.05	+0.03	+0.01	+0.02
SPLADE+FlanT5Base (0.470)	-0.01	-0.01	+0.01	0.00	+0.03
SPLADE+FlanT5XXL (0.495)	-0.02	+0.14	+0.09	+0.03	+0.01
Contriever+FlanT5Small (0.377)	-0.03	+0.07	0.00	0.00	+0.03
Contriever+FlanT5Base (0.478)	-0.02	+0.03	+0.06	-0.02	+0.03
Contriever+FlanT5XXL (0.496)	-0.02	+0.18	+0.04	+0.02	+0.04

Table 10: LaMP-5

		Fa	irness Interv	als	
Model (baseline utility)	[0.0, 0.2)	[0.2, 0.4)	[0.4, 0.6)	[0.6, 0.8)	[0.8, 1.0)
BM25+FlanT5Small (0.425)	-0.12	-0.07	-0.04	-0.07	-0.03
BM25+FlanT5Base (0.421)	-0.09	-0.03	-0.03	-0.04	-0.03
BM25+FlanT5XXL (0.536)	-0.03	+0.03	+0.04	-0.01	+0.05
SPLADE+FlanT5Small (0.362)	-0.06	-0.02	+0.03	+0.02	+0.01
SPLADE+FlanT5Base (0.361)	-0.02	+0.02	+0.03	+0.04	+0.05
SPLADE+FlanT5XXL (0.527)	-0.02	+0.03	+0.03	+0.03	+0.05
Contriever+FlanT5Small (0.351)	-0.05	-0.02	+0.03	+0.01	+0.04
Contriever+FlanT5Base (0.373)	-0.04	+0.01	0.00	+0.04	+0.03
Contriever+FlanT5XXL (0.526)	-0.02	+0.01	+0.02	+0.02	+0.06

Table 11: LaMP-6

	Fairness Intervals				
Model (baseline utility)	[0.0, 0.2)	[0.2, 0.4)	[0.4, 0.6)	[0.6, 0.8)	[0.8, 1.0]
BM25+FlanT5Small (0.490)	-0.02	-0.01	+0.01	+0.03	0.00
BM25+FlanT5Base (0.673)	-0.04	-0.04	-0.11	-0.03	-0.05
BM25+FlanT5XXL (0.626)	+0.02	-0.05	-0.05	-0.04	+0.03
SPLADE+FlanT5Small (0.525)	-0.05	-0.07	-0.04	-0.01	-0.03
SPLADE+FlanT5Base (0.659)	-0.03	-0.06	-0.06	-0.10	-0.02
SPLADE+FlanT5XXL (0.518)	+0.07	+0.06	+0.09	+0.04	+0.05
Contriever+FlanT5Small (0.440)	+0.02	+0.03	+0.05	+0.06	+0.08
Contriever+FlanT5Base (0.580)	+0.06	+0.04	0.00	+0.02	+0.06
Contriever+FlanT5XXL (0.607)	-0.01	-0.02	-0.03	-0.04	-0.04

Table 12: LaMP-7

J DATA STATISTICS

J.1 LAMP DATA STATISTICS FOR FLAN-T5-SMALL

1570	Dataset	#queries	Avg # Docs (Std)	Avg # Pos Labels	Avg % Pos
1571		•		(Std)	Labels
1572	LaMP-1	51	123.51 (82.66)	9.08 (11.63)	9.53
1573	LaMP-2	192	52.81 (46.21)	7.98 (9.64)	22.53
1574	LaMP-3	311	189.82 (134.33)	65.88 (95.77)	34.28
1575	LaMP-4	833	192.19 (195.28)	40.72 (61.82)	27.1
1576	LaMP-5	826	106.06 (71.47)	26.18 (31.1)	24.83
1577	LaMP-6	760	86.0 (52.66)	27.78 (29.22)	35.92
1578	LaMP-7	365	19.36 (18.4)	8.23 (10.38)	45.48

Table 13: LaMP data statistics for Flan-T5-Small after filtering for fairness evaluation.

J.2 LAMP DATA STATISTICS FOR FLAN-T5-BASE

1584	Dataset	#aueries	Avg # Docs (Std)	Avg # Pos Labels	Avg % Pos
1585		1	8	(Std)	Labels
1586	LaMP-1	232	102.86 (61.88)	20.07 (22.78)	22.49
1587	LaMP-2	280	45.58 (42.0)	10.45 (12.56)	29.87
1588	LaMP-3	378	185.32 (128.43)	73.82 (85.44)	41.19
1589	LaMP-4	827	186.98 (193.52)	49.79 (68.91)	31.57
1590	LaMP-5	759	105.62 (69.56)	26.09 (31.21)	25.71
1591	LaMP-6	783	86.18 (52.97)	30.11 (31.28)	38.65
1592	LaMP-7	211	21.72 (16.09)	6.96 (10.62)	33.02

Table 14: LaMP data statistics for Flan-T5-Base after filtering for fairness evaluation.

J.3 LAMP DATA STATISTICS FOR FLAN-T5-XXL

Dataset	#queries	Avg # Docs (Std)	Avg # Pos Labels	Avg % Pos
			(Std)	Labels
LaMP-1	264	111.66 (69.45)	25.12 (33.96)	23.35
LaMP-2	105	44.66 (42.82)	11.32 (15.61)	36.74
LaMP-3	182	198.06 (151.09)	41.86 (59.52)	22.19
LaMP-4	842	198.0 (200.82)	54.07 (73.34)	30.96
LaMP-5	511	104.18 (68.73)	23.1 (38.3)	23.39
LaMP-6	730	85.93 (52.46)	34.89 (35.54)	43.88
LaMP-7	151	20.6 (16.39)	8.58 (12.01)	42.7

Table 15: LaMP data statistics for Flan-T5-XXL after filtering for fairness evaluation.

Κ	IMPLEMENTATION DETAILS
BM	25:
	• Adapted from: https://github.com/dorianbrown/rank_bm25/tree/ master
SPI	ADE:
	 https://huggingface.co/naver/splade_v2_max
Con	triever:
	 https://huggingface.co/facebook/contriever
Flai	n-T5 Family:
	 https://huggingface.co/google/flan-t5-small
	 https://huggingface.co/google/flan-t5-base
	 https://huggingface.co/google/flan-t5-xxl
The	RAG model inferences were performed on NVIDIA A6000 GPUs with 48GB of VRAM.