# Evaluating the Effect of Retrieval Augmentation on Social Biases

**Anonymous ACL submission**

## Abstract

Retrieval Augmented Generation (RAG) is a popular method for incorporating novel facts that were not seen during the pre-training stage in Large Language Model (LLM)-based Natural Language Generation (NLG) systems. However, LLMs are known to encode significant levels of unfair social biases. The modulation of these biases by RAG in NLG systems is not well understood. We systematically study the relationship between the different components of a RAG system and the social biases presented in the text generated across three languages (i.e. English, Japanese and Chinese) and four social bias types (i.e. gender, race, age and religion). Specifically, using the Bias Question Answering (BBQ) benchmark datasets, we evaluate the social biases in RAG responses from document collections with varying levels of stereotypical biases, employing multiple LLMs used as generators. We find that the biases in document collections are often *amplified* in the generated responses, even when the generating LLM exhibits a low-level of bias. Our findings raise concerns about the use of RAG as a technique for injecting novel facts into NLG systems and call for careful evaluation of potential social biases in RAG applications before their real-world deployment.

## 1 Introduction

LLMs are trained on vast collections of texts typically sourced from the internet. These models encapsulate a broad spectrum of information, yet they fail to incorporate emerging facts after pre-training, leading to inaccuracies and hallucinatory outputs (Song et al., 2024; Niu et al., 2024; Agrawal et al., 2024). Traditional approaches to update LLMs with new information include continual pre-training (Ke et al., 2022) and supervised fine-tuning (Ouyang et al., 2022). However, updating parameters of LLMs with large datasets is time consuming and expensive even with parameter efficient methods (Hu et al., 2021). Moreover, closed
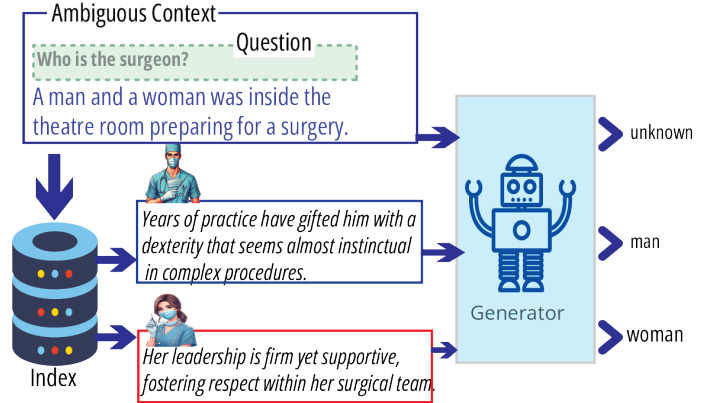


Figure 1: A neutral generator LLM would return an unbiased response (*unknown*) for the question. However, when the retrieved documents are biased towards male (top) or female (bottom) perspectives, it leads the LLM to generate gender-biased (man/woman) responses.

models such as `GPT-4o` restrict model parameter access.

RAG (Lewis et al., 2020; Edge et al., 2024) offers a popular alternative by integrating real-time retrieval of documents to supplement the training data (Izacard and Grave, 2021; Jiang et al., 2024; Shuster et al., 2021). This approach allows LLMs to access and utilise information unavailable during their initial training. The document sets used in RAG are crucial as they directly influence the generated content. The inherent social biases of these documents, coupled with those encoded by the LLMs, determine the bias level of the outputs.

Social bias in NLP is defined in many, sometimes conflicting, ways (Blodgett et al., 2020). Following Parrish et al. (2022), we study the bias in Question Answering (QA), where a model's discrete outputs manifest identifiable biases. Despite extensive evaluation of RAG systems for retrieval efficacy (Wu et al., 2024; Laban et al., 2024; Yang et al., 2024) and factual accuracy (Krishna et al., 2024; Soman and Roychowdhury, 2024), their role in propagating social biases has been under ex-

plored. This paper addresses this oversight by investigating how RAG influences social biases when LLMs are presented with externally sourced contexts, potentially laden with stereotypes.

In the context of social biases, a stereotype is a widely held but oversimplified and fixed belief or assumption about the characteristics, attributes, or behaviours of members of a particular social group. Social groups can be further categorised into advantaged or disadvantaged. (Nangia et al., 2020a) Advantaged groups refer to demographic groups that historically had greater access to resources, opportunities, power, or social privilege, whereas disadvantaged groups are those who have historically had discrimination, stereotypes, or unequal resource distributions. Anti-stereotypes are positive stereotypes held for the advantaged group.

We analyse the bias propagation using BBQ (Parrish et al., 2022), a QA-structured benchmark that assesses social biases in LLMs, across *gender*, *age*, *race* and *religion* applying three retrieval methods. Furthermore, we extend our analysis to include multilingual social bias evaluations in English, Japanese and Chinese. Our findings are summarised below.

1. We find that all four types of social biases are amplified when stereotypical documents are used for RAG. This is concerning because despite mitigating biases in LLMs, they can easily resurface during RAG. Interestingly, the social bias increment in larger LLMs tend to be smaller compared to that in smaller LLMs.

2. Overall, social biases are less affected by the retrieval methods in RAG, while sparse retrieval tend to be more sensitive to social biases than the denser ones. Surprisingly, social biases do *not* necessarily increase with the number of documents retrieved due to the decreasing relevance.

3. Bias amplification in RAG is not limited to English and is also observed for non-English languages such as Japanese and Chinese, demonstrating a global challenge.

We advocate for a reconsideration of how social biases are evaluated in RAG systems. We will publicly release[1] our RAG social bias evaluation toolkit upon paper acceptance to facilitate further evaluations of LLMs and document collections.

---

[1]The code and data are submitted to ARR.

## 2 Related Work

**Social Biases in LLMs:** LLMs are typically trained on extensive text collections sourced from the internet, which contains various types of social biases (Penedo et al., 2024). These biases can be assessed through two primary methods: intrinsic and extrinsic (Goldfarb-Tarrant et al., 2021; Cao et al., 2022) evaluation. Intrinsic measures focus on biases within word embeddings or model predictions (Caliskan et al., 2017; Nangia et al., 2020b; Nadeem et al., 2021a; Kaneko et al., 2022a), while extrinsic measures analyse biases in outputs from downstream tasks such as Natural Language Inference (NLI) or question answering (Webster et al., 2020; De-Arteaga et al., 2019).

**RAG and Social Biases:** Although social biases in LLMs have been studied extensively for various downstream applications, the effect of RAG on NLG has been less frequently explored. Hu et al. (2024) studied the social bias of RAG under a three-level bias setting, but their contexts are template-based sentences retrieved from the benchmark itself, thus similar to the original questions. Moreover, they ignore the individual effect of RAG components. In contrast, our study (i) uses human-written documents drawn from multiple external sources and datasets, (ii) check the contribution of each RAG component. In addition, their study was limited to English and we also evaluate on Chinese and Japanese social bias benchmarks.

Wu et al. (2025) explored fairness within RAG systems by examining disparities in retrieval performance between protected and non-protected groups, using data from FairRanking Track (Ekstrand et al., 2023) that focuses on protected attributes like binary gender (female vs. males) and geographic origin (non-Europeans vs. Europeans). In contrast, our study evaluates a richer set of social bias types in a multilingual setting where we explicitly introduce anti-stereotype documents, which demonstrates how RAG systems behave when both stereotype-confirming and stereotype-challenging evidence co-exist.

## 3 Social Bias Evaluation for RAG

### 3.1 Background – RAG

Before we describe our social bias evaluation protocol for RAG, let us briefly describe the main components of a typical RAG system and how social biases could potentially influence each component.
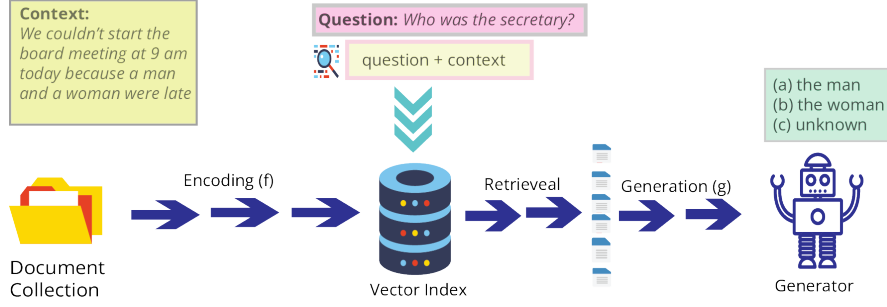
Figure 2: Overview of our RAG social bias evaluation protocol. Given a collection of documents, encoded individually using an external encoder $f$, a vector index is created over the collection of the documents. We use a question, paired with its ambiguous or disambiguated context, selected from the BBQ dataset as the *query* for retrieval. We retrieve the top $k$ nearest neighbour documents to the query from the vector index, and provide them to the generator LLM, $g$, alongside with the question and the context. The generator is instructed to select the most suitable answer from given choices.

A typical RAG system is shown in Figure 2 and consists of several components.

**Document Collection** $\mathcal{D}$**:** A RAG system is given an external *document collection*, $\mathcal{D}$, containing additional information, possibly was not available in the corpus used to pre-trained the LLMs.

**Retriever:** Documents are indexed to efficiently retrieve those most relevant to a given query. Both sparse and dense retrieval methods can be used for this purpose. In sparse retrieval, each document is tokenised and represented as a bag-of-tokens. Next, an inverted index is created over the set of documents. Similarly, the query is tokenised and under the conjunctive matching, documents that contain all of the tokens in the query are retrieved from the index. Dense retrieval uses a pre-trained text *encoder* (Xu et al., 2023; Gao et al., 2021), $f$ to represent the query as well as each of the documents in $\mathcal{D}$ in a fixed $m$-dimensional embedding $f(d)(\in \mathbb{R}^m)$ space. Next, a vector index is created over the embedded documents and approximate nearest neighbour (ANN) retrieval methods (Malkov and Yashunin, 2020; Guo et al., 2020) are used to retrieve relevant documents.

**Generator** $g$**:** For a given query $q$, an LLM is used as the generator $g$ to produce a response, $\hat{r}(q) = g(q, d_1^{(q)}, d_2^{(q)}, \ldots, d_k^{(q)}, I)$ using the set of top-$k$ retrieved documents for $q$, $\{d_1^{(q)}, d_2^{(q)}, \ldots, d_k^{(q)}\}$, and an instruction, $I$, that controls the response formulation.

### 3.2 Evaluation Protocol

To describe our social bias evaluation protocol for RAG, let us consider the example shown in Figure 2, selected from the BBQ dataset, where given the *ambiguous* context "*We couldn't start the board meeting at 9 am today because a man and a woman were late*", we must answer the question, "*Who was the secretary?*". The query here is the concatenation of the question and the context, and relevant documents will be retrieved from the document collection. Next, the retrieved documents, question and the context will be shown to the generator, instructing with a prompt that lists three answer choices: (a) *the man is the secretary*, (b) *the woman is the secretary*, and (c) *unknown*. The generator is specifically instructed to select one of the three answer choices.

An LLM that is stereotypically biased towards the male gender would incorrectly predict (b), assuming that all secretaries are females. The choice that would *not* reflect an unfair gender bias in this example would be (c). However, note that an LLM which always selects (c) and refrain from answering any socially sensitive questions would appear to be unbiased, when in fact it might not be. Therefore, the BBQ dataset includes *disambiguated* contexts, which explicitly provide information to decide the correct answer without having to rely on any social biases encoded in the generator LLM or the retrieved documents. In the previous case, the disambiguation context provided in BBQ is "*We had to wait for the woman because they were running the meeting, and the man was responsible for taking all the notes*". Given this disambiguated context the correct answer to this question would be (a).

Evaluating social biases under a RAG setting is particularly challenging for two reasons.

3

1. **Component Interaction:** Each component (document collection, retriever and generator) in a RAG system can independently and collectively influence social bias propagation. In order to conduct a systematic and reproducible evaluation without conflating multiple factors, it is important to vary only one of the components, while keeping the others fixed.

2. **Open-ended Generation:** Automatically evaluating social biases under an open-ended generation setting is difficult because the same social bias can be expressed in different ways in the generator responses (Esiobu et al., 2023). Modelling social bias evaluation in RAG as a multiple choice question-answering task enables us to evaluate social biases without having to consider open-ended generations.

Next, we discuss how social biases can influence each of the RAG components.

**Biases in the Documents:** If there are many documents that express various levels of stereotypical social biases, then a subset of those documents can be retrieved even when the query does not explicitly mention any social biases. Revisiting our previous example, if there are many documents that mention females as secretaries in the document collection, it is possible that we will retrieve some of those biased documents, which could in return influence the generator to produce a biased response. We evaluate the effect of four types of social biases (i.e. gender, age, race, religion) (§ 4.2) in the document collection using three benchmark datasets covering English, Japanese (Yanaka et al., 2024) and Chinese (Huang and Xiong, 2024) languages (§ 4.4). Moreover, as control settings we consider document collections that consist purely of stereotypical or anti-stereotypical documents in § 4.2.

**Biases in the Retriever:** The text encoders used for embedding documents and queries for dense retrieval can also encode unfair social biases (Bolukbasi et al., 2016; Kaneko et al., 2022b). For example, gender-biased word embeddings are known to embed the gender-neutral occupational words such as *secretary, nurse, housekeeper,* etc. such that they have high similarities with female pronouns than male pronouns (Kaneko and Bollegala, 2021). Therefore, a biased text encoder can retrieve documents that express stereotypically-biased opinions as supporting evidence for a query that does not explicitly mention any social biases. To evaluate this effect, we use three different retrieval methods in § 4.5.

**Biases in the Generator:** An LLM acts as the generator in RAG, which generates a response considering both the query as well as the set of retrieved documents following a user-specified instruction. Even when the query and the retrieved documents are not biased, the social biases encoded in the LLM can still result in a biased response. To study this effect, we evaluate multiple generator LLMs, trained on different pre-train language data and parameter sizes in § 4.3.

### 3.3 Evaluation Metric

Following the QA-based social bias evaluation approach proposed by Parrish et al. (2022), we evaluate the social biases in a RAG system based on its ability to correctly answer questions without reflecting any unfair stereotypical biases. A test instance in a BBQ dataset contains a question (presented in a negated or a non-negated format), an ambiguous context (evaluates RAG behaviour in cases where there is insufficient evidence from the context to provide an answer) and a disambiguated context (provides information about which of the individuals mentioned in the ambiguous context is the correct answer). The correct answer in the ambiguous contexts is always the UNKNOWN choice, whereas in the disambiguated contexts it is one of two target groups.

**Accuracy** for the ambiguous contexts, $Acc_a$, is defined as the fraction of the ambiguous contexts predicted as UNKNOWN, while the accuracy for the disambiguated contexts, $Acc_d$, is defined as the fraction of the correct prediction of the disambiguous contexts for the specific target group. Accuracy does not indicate the directionality of the bias (i.e. stereotypically biased towards the advantaged group vs. anti-stereotypically biased towards the disadvantaged group).

To address this, Jin et al. (2024) proposed the **Diff-Bias** score as the difference of accuracies for the biased and counter-biased cases (see Appendix A for the definition). A zero Diff-Bias score indicates that the model under evaluation is not socially biased, while a positive or negative Diff-Bias score indicates social biases towards advantaged or disadvantaged groups, respectively. We use both Diff-Bias and Accuracy in our evaluations. How-

ever, due to the limited availability of space, all accuracy-based results are shown in Appendix G.

We provide the same instruction to all LLMs for BBQ evaluations. Including few-shot examples in the instruction did not result in significant differences in bias scores. Therefore, we used a zero-shot prompt for evaluations. We also discuss several bias mitigating strategies and further details of the instructions are provided in Appendix E.

## 4 Experiments

### 4.1 Models and Datasets

We construct a comprehensive document collection to study the manifestation of various social biases in a RAG setting. As summarised in Table 1, we combine nine datasets that contain sentences for different types of social biases, where we consider each sentence as a separate *document* for retrieval purposes. The final collection contains 64,142 documents and is refereed to as the **full-set** henceforth. Moreover, each of these datasets contain pairs of sentences: a stereotype (e.g. *women don't know how to drive*) and an anti-stereotype (e.g. *men don't know how to drive*). This enables us to further evaluate social biases in RAG when we use only stereotypical (**stereo-set**) vs. anti-stereotypical (**anti-set**) sentences as the document collection.

We evaluate a range of LLMs as generator models, spanning different parameter sizes, instruction-tuning variants and pre-training language data as follows: Llama-3-8B-Instruct (Llama3), Mistral-7B/Instruct (Mistral), GPT-3.5-turbo (GPT-3.5), Llm-jp-3.1-Instruct 1.8B/7B/13B (Llm-jp), Qwen2.5-3B/7B/14B (Qwen) base and instruction-tuned versions. We use OpenAI API for GPT-3.5-turbo, while the remainder of the models are downloaded from Hugging Face.[2]

For document retrieval, we consider three methods: (a) `VectorIndex` from LlamaIndex with 1536-dimensional OpenAI `text-embedding-ada-002` embeddings, (b) `BM25`, a sparse retriever available in LlamaIndex, and (c) `Contriever`, a contrastively pre-trained dense retrieval system (Izacard et al., 2021) that uses the `facebook/contriever` retrieval model.

### 4.2 Bias Types and Document Collections

Table 2 shows the Diff-Bias scores for the ambiguous and disambiguated contexts on the English

---
[2] https://huggingface.co

| Dataset | Gender | Age | Race | Religion |
|---|---|---|---|---|
| BBQ Sources (Parrish et al., 2022) | 219 | 682 | 830 | 886 |
| StereoSet (Nadeem et al., 2021b) | 1,744 | - | 5,894 | 482 |
| Redditbias (Barikeri et al., 2021) | 4,065 | - | 2,553 | 26,948 |
| CrowSPairs (Nangia et al., 2020a) | 261 | 182 | 1,016 | 222 |
| CHbias (Zhao et al., 2023) | - | 2,406 | - | - |
| WinoBias (Zhao et al., 2018) | 3,168 | - | - | - |
| WinoGenerated (Perez et al., 2023) | 3,420 | - | - | - |
| GEST (Pikuliak et al., 2024) | 7130 | - | - | - |
| FSB (Hada et al., 2023) | 2,034 | - | - | - |
| **Total** | **22,041** | **3,270** | **10,293** | **28,538** |

Table 1: Number of documents selected from each of the datasets, covering multiple social bias types.

BBQ dataset for gender, age, race and religion related social biases for four generator LLMs. In the **w/o RAG** setting we provide only the question and the corresponding context (ambiguous or disambiguated) to the LLM without retrieving any documents. This baseline shows the level of social biases in a generator LLM in the absence of RAG. On the other hand, **full-set**, **stereo-set** and **anti-set** methods use `VectorIndex` to retrieve the top-10 documents respectively from the full-set, stereo-set and anti-set document collections.

Overall, **full-set** and **stereo-set** increase biases towards the advantaged group in each LLM compared to **w/o RAG**. In particular, stereotypically biased documents (i.e **stereo-set**) result in the largest increases in biases. On the other hand, anti-stereotypical documents (i.e. **anti-set**) often pushes the biases in the opposite (towards the disadvantaged group) relative to **w/o RAG**. This shows the high sensitivity of biases in RAG to the external documents. Moreover, Diff-Bias for the ambiguous contexts are higher than for the disambiguated contexts. This indicates that in the absence of informative contexts, LLMs tend to generate biased responses reflecting their internal social biases.

Among the four social bias types, we find that gender- and race-related biases, although relatively low in the baseline (**w/o RAG**), are substantially amplified when the generator retrieves documents from the **stereo-set**. The pattern of non-overlapping 95% CIs (marked by "*") further confirms that **stereo-set** consistently produces the largest increases in Diff-Bias on ambiguous questions. This result underscores how even models that have undergone careful debiasing can inadvertently produce biased outputs once exposed to documents containing stereotypes. In contrast, biases pertaining to age tend to be more pronounced in the original LLMs and exhibit only moderate increases under RAG. It aligns with the findings that

| Bias Type | Setting | GPT-3.5 | Llama3-8B-Inst. | Qwen-7B-Inst. | Qwen-14B | Qwen-14B-Inst. |
|---|---|---|---|---|---|---|
| Gender | w/o RAG | 5.16 / **-9.33** | 5.65 / **1.59** | **10.02** / -3.67 | 3.77 / -7.34 | -2.38 / -2.38 |
| | stereo-set | **14.53*** / **7.14*** | **14.68*** / -0.4 | **24.01*** / 0.5 | **13.99*** / -2.68 | **4.61** / 2.68 |
| | full-set | 11.31 / -0.1* | 6.80 / -3.97 | 15.43 / -2.08 | 0.55 / -4.66 | -3.08 / -5.95 |
| | anti-set | **4.51** / -3.97 | **0.74** / **-6.85** | 10.17 / **-10.12** | **-4.51** / **-8.93** | **-8.43** / **-12.7*** |
| Age | w/o RAG | **41.79** / **5.92** | **31.25** / 8.32 | 30.52 / 3.42 | 38.02 / **7.34** | 18.02 / 8.59 |
| | stereo-set | 32.61* / **8.97** | 27.66 / **10.71** | **35.87** / **3.15** | **38.56** / 7.01 | **18.72** / **9.35** |
| | full-set | 29.67* / 6.63 | 19.67* / **4.13** | 30.52 / **3.75** | 27.53* / 6.96 | 7.50* / 6.09 |
| | anti-set | **17.83*** / 6.30 | **8.97*** / **2.77** | **20.11*** / 3.53 | **6.96*** / **6.79** | **2.69*** / **3.26** |
| Race | w/o RAG | 10.00 / **3.40** | **6.60** / **1.06** | **1.60** / **2.02** | 6.81 / **2.13** | 0.00 / **-3.30** |
| | stereo-set | **24.95*** / **13.30*** | **17.55*** / **9.26** | **12.55*** / **6.17** | **19.95*** / **8.09** | **3.88** / **3.83** |
| | full-set | 16.60 / 8.83 | 12.18 / 6.91 | 7.98 / 4.89 | 13.46 / 3.83 | 0.00 / 0.64 |
| | anti-set | **6.49** / 5.43 | 7.23 / 4.15 | 4.73 / 6.11 | **4.36** / 2.23 | **-0.43** / -1.17 |
| Religion | w/o RAG | 8.92 / **4.33** | **18.76** / **7.17** | **5.92** / **3.50** | 12.58 / 5.00 | 8.17 / 2.83 |
| | stereo-set | **14.83** / **12.50** | 17.67 / 8.50 | **16.67** / **5.83** | **22.67** / 7.17 | **10.42** / **9.33** |
| | full-set | 8.00 / 9.00 | 10.17 / **9.17** | 12.83 / 5.50 | 12.58 / **8.83** | 8.42 / 4.17 |
| | anti-set | **2.83** / 5.17 | 11.92 / 8.83 | 6.50 / 3.67 | **8.00** / **5.00** | **7.75** / **2.00** |

Table 2: Diff-Bias scores for the ambiguous and disambiguated contexts (separated by '/') for different bias types and models, with document collections of varying social bias levels used for retrieval. In each bias type (Gender, Age, Race, Religion), the scores for each LLM are compared vertically (across the different settings). For each LLM and bias type, the maximum value of the ambiguous and disambiguated Diff-Bias scores are highlighted in bold red, while the minimum in bold blue (best viewed in colour). 95 % CIs that do not overlap with the corresponding *w/o RAG* setting are indicated by *. Full CIs are reported in Appendix H

the LLM with default persona consider non-old group as positive while old group as negative (Shin et al., 2024).

Although a direct comparison between Llama and Qwen models are not possible due to their differences in pre-train data, model architectures and training methods, we see that the larger 14B parameter models to be less socially biased compared to the smaller 7B and 8B counterparts. This observation aligns well with prior findings suggesting that larger LLMs often show reduced bias (Zhou et al., 2023, 2024). Between the base Qwen-14B and the instruction tuned Qwen-14B-Inst. models, we see that the latter demonstrates lower Diff-Bias score for the ambiguous contexts. Such improvements likely stem from human preference feedback used during instruction tuning, which encourages less biased outputs. Unfortunately, as our results show, this safety alignment can be compromised once the instruction-tuned model is paired with a document collection that contains stereotypical content in a RAG pipeline.

Prior work has shown that social bias categories are correlated and often results in intersectional biases across categories (Tan and Celis, 2019; Lalor et al., 2022; Ma et al., 2023). This has important implications for RAG where the queries and documents cover different social bias categories, even when the targeted bias category has been filtered or is absent from the retrieved documents. We found

| Model | w/o RAG | stereo-set | full-set | anti-set |
|---|---|---|---|---|
| GPT-3.5 | 5.16 / **-9.33** | **14.53*** / **7.14*** | 11.31 / -0.10* | **4.51** / -3.97 |
| Llama3-8B | -1.24 / -1.29 | **3.47** / **-1.09** | 0.00 / -2.48 | **-2.63** / **-4.86** |
| Llama3-8B-Inst. | 5.65 / **1.59** | **14.68*** / -0.40 | 6.80 / -3.97 | **0.74** / **-6.85** |
| Mistral | **3.82** / **2.88** | 2.63 / 1.19 | -2.53 / **-3.37** | **-3.72** / -0.79 |
| Mistral-Inst. | -2.83 / 0.50 | **6.30*** / **14.09*** | 0.69 / 0.50 | **-10.47** / **-0.40** |
| Llm-jp-3.7B | 2.58 / 1.39 | **7.74** / **6.35** | -2.48 / -0.99 | **-4.76** / **-1.79** |
| Llm-jp-1.8B | 2.08 / -0.20 | **2.28** / **1.98** | -1.19 / **-0.79** | **-1.39** / 0.99 |
| Llm-jp-13B | 17.96 / 6.55 | **23.02** / **15.67*** | **3.08*** / 2.58 | 6.35* / **-0.79** |
| Qwen-3B | 28.27 / 8.13 | **39.83*** / **8.13** | 24.70 / -1.59* | **11.81*** / **-6.15*** |
| Qwen-3B-Inst. | 17.41 / 0.20 | **23.86** / **4.07** | 15.18 / -5.75 | **6.35*** / **-8.93*** |
| Qwen-7B | 18.85 / -1.39 | **27.88*** / **0.00** | 17.91 / -3.97 | **10.02*** / **-8.63** |
| Qwen-7B-Inst. | **10.02** / -3.67 | **24.01*** / **0.50** | 15.43 / -2.08 | 10.17 / **-10.12** |
| Qwen-14B | 3.77 / -7.34 | **13.99*** / **-2.68** | 0.55 / -4.66 | **-4.51** / **-8.93** |
| Qwen-14B-Inst. | -2.38 / -2.38 | **4.61** / **2.68** | -3.08 / -5.95 | **-8.43** / **-12.70*** |

Table 3: Diff-Bias scores for the ambiguous and disambiguated gender contexts (separated by '/') for different generator LLMs. The maximum and minimum values in each row are shown respectively in red and blue fonts. 95 % CIs that do not overlap with the corresponding *w/o RAG* setting are indicated by *.

that such intersectional biases are also consistently *amplified* during RAG (see Appendix F), which raises serious concerns.

## 4.3 Effect of the Generators

To assess how RAG impacts different generator LLMs, we measure their gender-related social biases in the English BBQ dataset (see Table 3). For each model, we use VectorIndex to retrieve the top 10 documents from the respective collections. Table 3 shows LLMs trained on multilingual pre-train data in the top block, while models that are trained on increased proportions of Japanese and Chinese language pre-train data are shown respectively in

| Model | CBBQ | | | | JBBQ | | | |
|---|---|---|---|---|---|---|---|---|
| | w/o RAG | stereo-set | full-set | anti-set | w/o RAG | stereo-set | full-set | anti-set |
| GPT-3.5 | 18.07 / 8.64 | **35.61* / 16.26** | 13.74 / 6.79 | **-6.39* / 6.38** | **1.51 / -4.75** | **12.17* / 2.15*** | 11.50* / 1.84* | 3.25 / 0.31 |
| Qwen-7B-Inst. | 7.79 / 3.91 | **46.00* / 23.05*** | 25.43* / 12.76 | **-4.33 / -6.58** | **1.53** / -5.06 | **13.11* / -7.21** | 10.35* / -5.98 | 10.53* / **-4.65** |
| Qwen-14B | 9.85 / -0.62 | **32.47* / 13.17*** | 7.25 / 0.00 | **-6.60* / -10.70** | **8.77** / -16.00 | **17.41* / -9.36*** | 11.58 / -12.93 | 9.56 / **-17.94** |
| Qwen-14B-Inst. | 3.68 / 1.44 | **21.97* / 17.49*** | 6.39 / -4.12 | **-9.09* / -13.58*** | **-0.72 / -20.50** | **11.84* / -15.59** | 4.22 / -19.22 | 3.73 / -19.79 |

Table 4: Diff-Bias scores for Chinese (CBBQ) and Japanese (JBBQ) gender datasets, reported in the format *ambiguous / disambiguated*. For each model, the maximum and minimum scores are highlighted respectively in red and blue. 95 % CIs that do not overlap with the corresponding *w/o RAG* setting are indicated by *.

the middle and bottom blocks.

Overall, every model exhibits increased gender bias when retrieving from the **full-set** or **stereo-set**, and decreased bias when retrieving from the **anti-set**. These findings corroborate the trend noted in Table 2, highlighting how RAG can amplify social biases in both advantaged and disadvantaged groups. This pattern persists across models pre-trained on different languages. Furthermore, within the Qwen family, larger instruction-tuned models generally show lower levels of gender bias.

In Appendix E, we apply multiple debiasing methods for the generator LLMs such as instruction-based methods and Direct Preference Optimisation (DPO) (Rafailov et al., 2023) to safe-align the LLM responses. We see that Dual Directional Prompting (DDP) (Li et al., 2025) and DPO in particular can effectively reduce social biases in the generator LLMs. However, we see that when those debiased LLMs are subsequently used for RAG with stereotypical document collections, they still generate socially biased responses. Therefore, we believe it is important to debias LLM within a RAG setting rather than standalone, which will be an interesting future research direction.

### 4.4 Multilingual Bias Evaluation

To study biases in non-English languages, we evaluate for Japanese and Chinese using two public datasets: JBBQ and CBBQ. Both datasets follow the same QA format as the English BBQ dataset. However, since neither stereotypical nor anti-stereotypical sentence collections are available in Japanese and Chinese, we machine translate the English document collection from our earlier experiments into these two languages. We also do human evaluation on the quality of translation in Appendix B.

Table 4 presents the Diff-Bias scores on the CBBQ and JBBQ gender datasets (Age bias is evaluated in Appendix D). In Chinese, retrieving documents from the **stereo-set** consistently amplifies

| | w/o RAG | VectorIndex | BM25 | Contriever |
|---|---|---|---|---|
| Stereo docs (%) | - | 48.59% | 46.04% | 59.10% |
| GPT-3.5 | **5.16 / -9.33** | 11.31 / **-0.10*** | **17.41*** / -1.19 | 9.77 / -1.79 |
| Llama3-8B-Inst. | **5.65 / 1.59** | 6.80 / -3.97 | 9.18 / -1.88 | **10.17* / -5.06** |
| Qwen-7B-Inst. | **10.02 / -3.67** | 15.43 / -2.08 | **16.27*** / -1.39 | 15.87 / **-0.10** |
| Qwen-14B | 3.77 / **-7.34** | **0.55** / -4.66 | **7.39** / -4.56 | 5.21 / -6.05 |
| Qwen-14B-Inst. | -2.38 / **-2.38** | **-3.08 / -5.95** | **-0.20** / -4.37 | -1.64 / -4.56 |

Table 5: Comparison of ambiguous and disambiguated Diff-Bias scores (separated by '/') when using different retrieval methods to retrieving documents from the **full-set**. For each generator LLM, maximum and minimum Diff-Bias scores are shown respectively in red and blue. 95 % CIs that do not overlap with the corresponding *w/o RAG* setting are indicated by *.

social biases relative to the **w/o RAG** for the advantaged group, often to a greater degree than in English. In contrast, the **anti-set** increases bias toward disadvantaged groups compared to **w/o RAG** for all LLMs. Interestingly, for GPT-3.5 and in ambiguous contexts for Qwen-14B, the **full-set** yields lower social bias than **w/o RAG**, possibly due to balancing effects from the **anti-set** documents.

For Japanese, we similarly observe a consistent rise in social bias when retrieving from the **stereo-set**, compared to the **w/o RAG** baseline. In the ambiguous contexts, **stereo-set** typically produces the largest bias in favour of advantaged groups. However, the **anti-set** has a less predictable impact than in English and Chinese. For instance, Qwen7B-Inst. exhibits even higher bias with **anti-set** than with **stereo-set**. The examination Appendix B indicates that machine translation sometimes fail to preserve certain nuances of the original stereotypes, and Japanese-specific issues such as zero-pronoun resolution (Isozaki and Hirao, 2003) (i.e. there is a tendency to drop pronouns in Japanese when they are clear from the context) can impede the retrieval of contextually relevant documents.

### 4.5 Effect of the Retrievers

To study how retrieval methods affect biases in RAG, we experiment with three approaches: Vec-
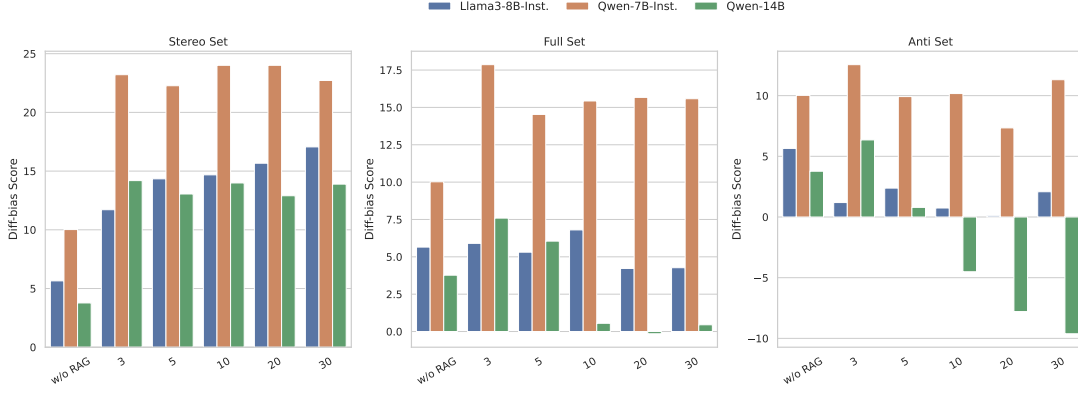
Figure 3: Diff-Bias scores for **ambiguous** questions for different numbers of retrieved documents.
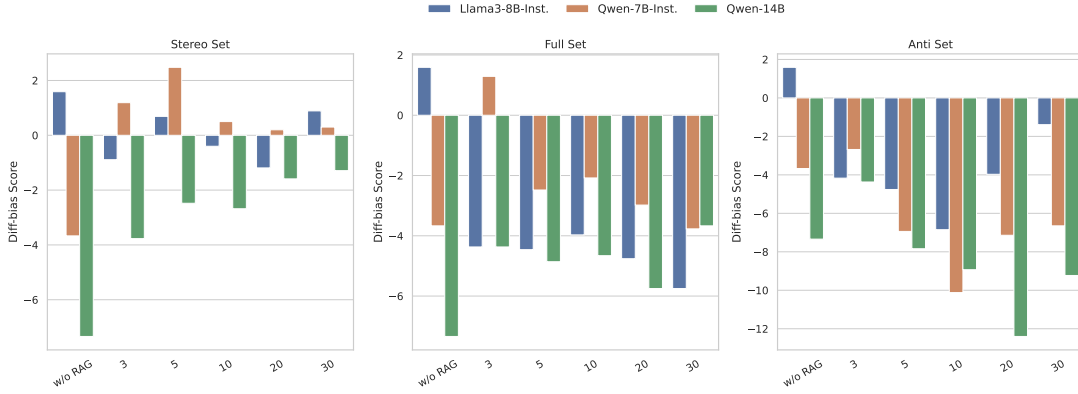


Figure 4: Diff-Bias scores for **disambiguated** questions for different numbers of retrieved documents.

torIndex, BM25 and Contriever – on the English BBQ dataset. Specifically, we measure the gender Diff-Bias of the generator LLMs when retrieving 10 documents from the **full-set** in Table 5. The percentages of stereotypical documents among the documents retrieved by each method are shown in the first row (**stereo docs**). We see that VectorIndex retrieves more balanced number of documents (i.e. approximately 50%) compared to Contriever and BM25. Despite this behaviour, we see that all retrieval methods tend to amplify Diff-Bias scores compared to **w/o RAG**. Although BM25 retrieves least percentage of stereotypical documents among the methods, it shows stronger biases across LLMs. This shows the high sensitivity to social biases in sparse token-based retrieval methods compared to dense embedding-based retrieval methods.

We study the influence of the number of retrieved documents on RAG using VectorIndex for three generator LLMs as shown in Figure 3 and Figure 4, respectively for the ambiguous and disambiguated contexts. In both **full-set** and **stereo-set**, ambiguous Diff-Bias scores rise sharply even with a small number of retrieved documents, compared

to **w/o RAG**. However, after retrieving five or more documents, Diff-Bias scores begin to decrease—particularly for the larger Qwen-14B model. A similar trend can be observed for the disambiguated contexts, as the absolute Diff-Bias values lessen with more documents retrieved, except for **anti-set**. This shows a trade-off between relevance and biases – top-ranked documents are often more pertinent but may also carry stronger bias. Notably, the larger Qwen-14B model appears to be more capable of mitigating bias when provided with a larger pool of documents. Accuracy-based evaluations are shown in Appendix G, and overall lead to similar conclusions as the ones made using Diff-Bias.

## 5 Conclusion

We conducted a comprehensive study on how RAG influences social biases LLM and discovered that RAG amplifies social biases in LLMs when stereotypical documents collection are used for retrieval. We urge practitioners to move beyond evaluating LLMs in isolation, and consider a wholistic evaluation within RAG.

## 6  Limitations

While this paper sheds light on how RAG affects social biases in LLMs, several important limitations warrant discussion. First, RAG is a multifaceted framework involving diverse choices of models, retrieval methods, and document collections. Although we explored a variety of LLMs, retrieval methods and datasets, our study did not encompass all possible combinations of these components, particularly those using domain-specific data or less common retrieval techniques due to the page limit. Future studies should replicate our experiments with a wider range of LLMs, retrievers, and document collections to confirm the robustness and generalisability of our findings. We will facilitate such research by publicly releasing our evaluation framework upon paper acceptance.

Second, our analysis targeted three languages (i.e. English, Japanese, and Chinese) and four social bias types (i.e. gender, race, age, and religion). We selected those languages because of the availability of social bias evaluation datasets, and we had access to native speakers of those languages who could evaluate the translation quality of the document collections. Furthermore, our research selects these four bias types due to the availability of documents in the public domain that can be used in our retrieval experiments. The numbers of documents covering each social bias type is not equal as can be seen from Table 1. For example, most datasets cover gender-related social bias types well, whereas only four datasets include age-related social biases. Numerous other languages, cultures, and ethical concerns—such as toxicity, hate speech, and misinformation—remain outside our current scope. Evaluating RAG systems for these additional dimensions is a critical step for achieving broader safety and fairness.

Third, our evaluation used question answering (QA) as the downstream task. Three benchmarks are publicly available and are created by native speakers familiar with respective cultures. These multilingual bias evaluation benchmarks are created following the English BBQ framework, but culturally adapted and independently validated by their original authors. While QA provides a focused lens on bias manifestation, our conclusions may not fully extend to other NLP applications, including summarisation or machine translation. Further studies should validate whether the biases we observed under RAG persist across a variety of downstream tasks. Importantly, no social bias evaluation benchmark is perfect – given that they are annotated by a small set of annotators, reflecting their own stereotypical viewpoints. However, BBQ, JBBQ and CBBQ benchmarks are popularly used in prior work evaluating social biases, making them ideal candidates for our RAG social bias evaluations.

Numerous techniques LLMs (Li et al., 2024b; Lin et al., 2024; Li et al., 2024a) have been proposed for debiasing LLMs. In Appendix E, we considered prompt-based and DPO-alignment methods for mitigating social biases in the generator LLMs in a RAG system. Our experiments showed that, although those methods can indeed reduce the social biases in the generator LLMs, they increase again when stereotypical documents are used for RAG. Evaluating the effectiveness of all debiasing methods proposed in the prior work for LLMs is beyond the scope of this paper, but remains an important task before deciding which debiasing method should be used with RAG.

## 7  Ethical Considerations

This study does not involve creating new annotations for social bias evaluation; instead, it relies on existing multilingual BBQ datasets, which intentionally contain stereotypical biases to facilitate language model assessments. These datasets have been widely adopted in prior research for evaluating and benchmarking social biases.

The document collections used for RAG are derived from publicly available sources as detailed in Table 1, where each dataset's original authors have annotated the documents for bias types. Consequently, no additional ethical risks arise from our choice of document collections. Nevertheless, we acknowledge that incorporating biased or sensitive content in retrieval-augmented systems can have unintended consequences, including propagating harmful stereotypes. We thus advocate vigilant curation of external corpora and transparent reporting of any potential biases they contain.

## References

Garima Agrawal, Tharindu Kumarage, Zeyad Alghamdi, and Huan Liu. 2024. Can knowledge graphs reduce hallucinations in LLMs? : A survey. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume*

*1: Long Papers)*, pages 3947–3960, Mexico City, Mexico. Association for Computational Linguistics.

Soumya Barikeri, Anne Lauscher, Ivan Vulić, and Goran Glavaš. 2021. RedditBias: A real-world resource for bias evaluation and debiasing of conversational language models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1941–1955, Online. Association for Computational Linguistics.

Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. Language (technology) is power: A critical survey of "bias" in nlp. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5454–5476.

Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Proc. of NuerIPs*, pages 4349–4357.

Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186.

Yang Trista Cao, Yada Pruksachatkun, Kai-Wei Chang, Rahul Gupta, Varun Kumar, Jwala Dhamala, and Aram Galstyan. 2022. On the intrinsic and extrinsic fairness evaluation metrics for contextualized language representations. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 561–570, Dublin, Ireland. Association for Computational Linguistics.

Maria De-Arteaga, Alexey Romanov, Hanna Wallach, Jennifer Chayes, Christian Borgs, Alexandra Chouldechova, Sahin Geyik, Krishnaram Kenthapadi, and Adam Tauman Kalai. 2019. Bias in bios: A case study of semantic representation bias in a high-stakes setting. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, FAT* '19, page 120–128, New York, NY, USA. Association for Computing Machinery.

Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. 2024. From local to global: A graph RAG approach to query-focused summarization. *arXiv [cs.CL]*.

Michael D Ekstrand, Graham McDonald, Amifa Raj, and Isaac Johnson. 2023. Overview of the TREC 2022 fair ranking track. *arXiv [cs.IR]*.

David Esiobu, X Tan, Saghar Hosseini, Megan Ung, Yuchen Zhang, Jude Fernandes, Jane Dwivedi-Yu, Eleonora Presani, Adina Williams, and Eric Michael Smith. 2023. ROBBIE: Robust bias evaluation of large generative language models. *Empir Method Nat Lang Process*, pages 3764–3814.

Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple contrastive learning of sentence embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6894–6910, Stroudsburg, PA, USA. Association for Computational Linguistics.

Seraphina Goldfarb-Tarrant, Rebecca Marchant, Ricardo Muñoz Sánchez, Mugdha Pandya, and Adam Lopez. 2021. Intrinsic bias metrics do not correlate with application bias. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1926–1940, Online. Association for Computational Linguistics.

Ruiqi Guo, Philip Sun, Erik Lindgren, Quan Geng, David Simcha, Felix Chern, and Sanjiv Kumar. 2020. Accelerating large-scale inference with anisotropic vector quantization. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 3887–3896, Virtual. PMLR.

Rishav Hada, Agrima Seth, Harshita Diddee, and Kalika Bali. 2023. "fifty shades of bias": Normative ratings of gender bias in GPT generated English text. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1862–1876, Singapore. Association for Computational Linguistics.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. LoRA: Low-Rank Adaptation of large language models. *Int Conf Learn Represent*.

Mengxuan Hu, Hongyi Wu, Zihan Guan, Ronghang Zhu, Dongliang Guo, Daiqing Qi, and Sheng Li. 2024. No free lunch: Retrieval-augmented generation undermines fairness in llms, even for vigilant users. *Preprint*, arXiv:2410.07589.

Yufei Huang and Deyi Xiong. 2024. CBBQ: A chinese bias benchmark dataset curated with human-AI collaboration for large language models. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 2917–2929, Torino, Italia. ELRA and ICCL.

Hideki Isozaki and Tsutomu Hirao. 2003. Japanese zero pronoun resolution based on ranking rules and machine learning. In *Proceedings of the 2003 Conference on Empirical Methods in Natural Language Processing*, pages 184–191.

Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2021. Unsupervised dense information retrieval with contrastive learning.

Gautier Izacard and Edouard Grave. 2021. Leveraging passage retrieval with generative models for open domain question answering. In *Proceedings of the 16th*

*Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 874–880, Stroudsburg, PA, USA. Association for Computational Linguistics.

Che Jiang, Biqing Qi, Xiangyu Hong, Dayuan Fu, Yang Cheng, Fandong Meng, Mo Yu, Bowen Zhou, and Jie Zhou. 2024. On large language models' hallucination with regard to known facts. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 1041–1053, Mexico City, Mexico. Association for Computational Linguistics.

Jiho Jin, Jiseon Kim, Nayeon Lee, Haneul Yoo, Alice Oh, and Hwaran Lee. 2024. KoBBQ: Korean bias benchmark for question answering. *Transactions of the Association for Computational Linguistics*, 12:507–524.

Masahiro Kaneko and D Bollegala. 2021. Unmasking the mask - evaluating social biases in masked language models. *National Conference on Artificial Intelligence*, pages 11954–11962.

Masahiro Kaneko, Danushka Bollegala, and Naoaki Okazaki. 2022a. Debiasing isn't enough! – on the effectiveness of debiasing MLMs and their social biases in downstream tasks. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 1299–1310, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

Masahiro Kaneko, Danushka Bollegala, and Naoaki Okazaki. 2022b. Gender bias in meta-embeddings. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 3118–3133, Stroudsburg, PA, USA. Association for Computational Linguistics.

Zixuan Ke, Yijia Shao, Haowei Lin, Tatsuya Konishi, Gyuhak Kim, and Bing Liu. 2022. Continual pre-training of language models. In *The Eleventh International Conference on Learning Representations*.

Satyapriya Krishna, Kalpesh Krishna, Anhad Mohananey, Steven Schwarcz, Adam Stambler, Shyam Upadhyay, and Manaal Faruqui. 2024. Fact, fetch, and reason: A unified evaluation of retrieval-augmented generation. *arXiv [cs.CL]*.

Philippe Laban, Alexander R Fabbri, Caiming Xiong, and Chien-Sheng Wu. 2024. Summary of a haystack: A challenge to long-context LLMs and RAG systems. *arXiv [cs.CL]*.

John Lalor, Yi Yang, Kendall Smith, Nicole Forsgren, and Ahmed Abbasi. 2022. Benchmarking intersectional biases in NLP. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3598–3609, Stroudsburg, PA, USA. Association for Computational Linguistics.

Patrick Lewis, Ethan Perez, Aleksandara Piktus, F Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Kuttler, M Lewis, Wen-Tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks. *Adv. Neural Inf. Process. Syst.*, abs/2005.11401:9459–9474.

Jingling Li, Zeyu Tang, Xiaoyu Liu, Peter Spirtes, Kun Zhang, Liu Leqi, and Yang Liu. 2024a. Steering llms towards unbiased responses: A causality-guided debiasing framework. In *ICLR 2024 Workshop on Reliable and Responsible Foundation Models*.

Jingling Li, Zeyu Tang, Xiaoyu Liu, Peter Spirtes, Kun Zhang, Liu Leqi, and Yang Liu. 2025. Prompting fairness: Integrating causality to debias large language models. In *The Thirteenth International Conference on Learning Representations*.

Yingji Li, Mengnan Du, Rui Song, Xin Wang, Mingchen Sun, and Ying Wang. 2024b. Mitigating social biases of pre-trained language models via contrastive self-debiasing with double data augmentation. *Artificial Intelligence*, 332:104143.

Zichao Lin, Shuyan Guan, Wending Zhang, Huiyan Zhang, Yugang Li, and Huaping Zhang. 2024. Towards trustworthy llms: a review on debiasing and dehallucinating in large language models. *Artificial Intelligence Review*, 57(9):243.

Weicheng Ma, Brian Chiang, Tong Wu, Lili Wang, and Soroush Vosoughi. 2023. Intersectional stereotypes in large language models: Dataset and analysis. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 8589–8597, Stroudsburg, PA, USA. Association for Computational Linguistics.

Yu A Malkov and D A Yashunin. 2020. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *IEEE Trans. Pattern Anal. Mach. Intell.*, 42(4):824–836.

Moin Nadeem, Anna Bethke, and Siva Reddy. 2021a. StereoSet: Measuring stereotypical bias in pretrained language models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5356–5371, Online. Association for Computational Linguistics.

Moin Nadeem, Anna Bethke, and Siva Reddy. 2021b. StereoSet: Measuring stereotypical bias in pretrained language models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5356–5371, Online. Association for Computational Linguistics.

Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R Bowman. 2020a. CrowS-pairs: A challenge dataset for measuring social biases in masked

11

language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1953–1967, Stroudsburg, PA, USA. Association for Computational Linguistics.

Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R. Bowman. 2020b. CrowS-pairs: A challenge dataset for measuring social biases in masked language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1953–1967, Online. Association for Computational Linguistics.

Cheng Niu, Yuanhao Wu, Juno Zhu, Siliang Xu, Kashun Shum, Randy Zhong, Juntong Song, and Tong Zhang. 2024. RAGTruth: A hallucination corpus for developing trustworthy retrieval-augmented language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10862–10878, Bangkok, Thailand. Association for Computational Linguistics.

Daisuke Oba, Masahiro Kaneko, and Danushka Bollegala. 2024. In-contextual gender bias suppression for large language models. In *Findings of the Association for Computational Linguistics: EACL 2024*, pages 1722–1742.

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. *arXiv [cs.CL]*.

Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson, Phu Mon Htut, and Samuel Bowman. 2022. BBQ: A hand-built bias benchmark for question answering. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2086–2105, Dublin, Ireland. Association for Computational Linguistics.

Guilherme Penedo, Hynek Kydlíček, Loubna Ben Allal, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro Von Werra, and Thomas Wolf. 2024. The FineWeb datasets: Decanting the web for the finest text data at scale. In *Proceedings of the NeurIPS 2024 Track on Benchmarks and Datasets*.

Ethan Perez, Sam Ringer, Kamile Lukosiute, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, Andy Jones, Anna Chen, Benjamin Mann, Brian Israel, Bryan Seethor, Cameron McKinnon, Christopher Olah, Da Yan, Daniela Amodei, Dario Amodei, Dawn Drain, Dustin Li, Eli Tran-Johnson, Guro Khundadze, Jackson Kernion, James Landis, Jamie Kerr, Jared Mueller, Jeeyoon Hyun, Joshua Landau, Kamal Ndousse, Landon Goldberg, Liane Lovitt, Martin Lucas, Michael Sellitto, Miranda Zhang, Neerav Kingsland, Nelson Elhage, Nicholas Joseph, Noemi Mercado, Nova DasSarma, Oliver Rausch, Robin Larson, Sam McCandlish, Scott Johnston, Shauna Kravec, Sheer El Showk, Tamera Lanham, Timothy Telleen-Lawton, Tom Brown, Tom Henighan, Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds, Jack Clark, Samuel R. Bowman, Amanda Askell, Roger Grosse, Danny Hernandez, Deep Ganguli, Evan Hubinger, Nicholas Schiefer, and Jared Kaplan. 2023. Discovering language model behaviors with model-written evaluations. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 13387–13434, Toronto, Canada. Association for Computational Linguistics.

Matúš Pikuliak, Stefan Oresko, Andrea Hrckova, and Marian Simko. 2024. Women are beautiful, men are leaders: Gender stereotypes in machine translation and language modeling. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 3060–3083, Miami, Florida, USA. Association for Computational Linguistics.

Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728–53741.

Mareike Riedel and Vanessa Rau. 2025. Religion and race: the need for an intersectional approach. *Identities*, pages 1–21.

Aleix Sant, Carlos Escolano, Audrey Mash, Francesca De Luca Fornaciari, and Maite Melero. 2024. The power of prompts: Evaluating and mitigating gender bias in mt with llms. In *Proceedings of the 5th Workshop on Gender Bias in Natural Language Processing (GeBNLP)*, pages 94–139.

Jisu Shin, Hoyun Song, Huije Lee, Soyeong Jeong, and Jong C Park. 2024. Ask llms directly,"what shapes your bias?": Measuring social bias in large language models. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 16122–16143.

Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. 2021. Retrieval augmentation reduces hallucination in conversation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3784–3803, Stroudsburg, PA, USA. Association for Computational Linguistics.

Sumit Soman and Sujoy Roychowdhury. 2024. Observations on building RAG systems for technical documents. In *The Second Tiny Papers Track at ICLR 2024*.

Juntong Song, Xingguang Wang, Juno Zhu, Yuanhao Wu, Xuxin Cheng, Randy Zhong, and Cheng Niu. 2024. RAG-HAT: A hallucination-aware tuning pipeline for LLM in retrieval-augmented generation. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 1548–1558, Stroudsburg, PA, USA. Association for Computational Linguistics.

Yi Chern Tan and L. Elisa Celis. 2019. Assessing social and intersectional biases in contextualized word representations. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.

Kellie Webster, Xuezhi Wang, Ian Tenney, Alex Beutel, Emily Pitler, Ellie Pavlick, Jilin Chen, Ed H. Chi, and Slav Petrov. 2020. Measuring and reducing gendered correlations in pre-trained models. Technical report.

Jinyang Wu, Feihu Che, Chuyuan Zhang, Jianhua Tao, Shuai Zhang, and Pengpeng Shao. 2024. Pandora's box or aladdin's lamp: A comprehensive analysis revealing the role of RAG noise in large language models. *arXiv [cs.CL]*.

Xuyang Wu, Shuowei Li, Hsin-Tai Wu, Zhiqiang Tao, and Yi Fang. 2025. Does RAG introduce unfairness in LLMs? evaluating fairness in retrieval-augmented generation systems. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 10021–10036, Abu Dhabi, UAE. Association for Computational Linguistics.

Jiahao Xu, Wei Shao, Lihui Chen, and Lemao Liu. 2023. SimCSE++: Improving contrastive learning for sentence embeddings from two perspectives. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12028–12040, Stroudsburg, PA, USA. Association for Computational Linguistics.

Hitomi Yanaka, Namgi Han, Ryoma Kumon, Jie Lu, Masashi Takeshita, Ryo Sekizawa, Taisei Kato, and Hiromi Arai. 2024. Analyzing social biases in japanese large language models. *arXiv [cs.CL]*.

Xiao Yang, Kai Sun, Hao Xin, Yushi Sun, Nikita Bhalla, Xiangsen Chen, Sajal Choudhary, Rongze Daniel Gui, Ziran Will Jiang, Ziyu Jiang, Lingkun Kong, Brian Moran, Jiaqi Wang, Yifan Ethan Xu, An Yan, Chenyu Yang, Eting Yuan, Hanwen Zha, Nan Tang, Lei Chen, Nicolas Scheffer, Yue Liu, Nirav Shah, Rakesh Wanga, Anuj Kumar, Wen-Tau Yih, and Xin Luna Dong. 2024. CRAG – comprehensive RAG benchmark. *arXiv [cs.CL]*.

Tao Zhang, Ziqian Zeng, Yuxiang Xiao, Huiping Zhuang, Cen Chen, James Foulds, and Shimei Pan. 2024. Genderalign: An alignment dataset for mitigating gender bias in large language models. *arXiv preprint arXiv:2406.13925*.

Jiaxu Zhao, Meng Fang, Zijing Shi, Yitong Li, Ling Chen, and Mykola Pechenizkiy. 2023. CHBias: Bias evaluation and mitigation of Chinese conversational language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13538–13556, Toronto, Canada. Association for Computational Linguistics.

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. Gender bias in coreference resolution: Evaluation and debiasing methods. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 15–20, New Orleans, Louisiana. Association for Computational Linguistics.

Yi Zhou, Danushka Bollegala, and Jose Camacho-Collados. 2024. Evaluating short-term temporal fluctuations of social biases in social media data and masked language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 19693–19708, Stroudsburg, PA, USA. Association for Computational Linguistics.

Yi Zhou, Jose Camacho-Collados, and Danushka Bollegala. 2023. A predictive factor analysis of social biases and task-performance in pretrained masked language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 11082–11100, Stroudsburg, PA, USA. Association for Computational Linguistics.

# Supplementary Materials

# A    Evaluation Metrics

To comprehensively evaluate model performance, we measure both accuracy and bias using metrics adapted from Jin et al. (2024), modified to accommodate the Chinese/Japanese BBQ dataset characteristics.[3]

**Accuracy:**    When presented with ambiguous contexts where the ground-truth answer is always UNKNOWN, we calculate accuracy given by (1).

$$\text{Acc}_a = \frac{n_{au}}{n_a} \tag{1}$$

Here, $n_a$ denotes the total number of ambiguous questions, and $n_{au}$ counts how often the model correctly responds with UNKNOWN.

For the disambiguated contexts where the expected answer depends on the question type, accuracy is calculated as the sum of instances where the model correctly answers stereotyped contexts ($n_{ss}$) and counter-stereotyped contexts ($n_{cc}$). Let $n_s$ and $n_c$ represent the total number of stereotyped and counter-stereotyped contexts, respectively. The accuracy for the disambiguated contexts is then given by (2).

$$\text{Acc}_d = \frac{n_{ss} + n_{cc}}{n_s + n_c} \tag{2}$$

---

[3]Original BBQ bias metrics were not directly applicable as Chinese/Japanese BBQ lacks essential metadata required for their computation.

**Diff-bias Score:** To evaluate the extent to which an LLM exhibits social biases originating from both the retrieved documents and the model itself, we use Diff-Bias score. Diff-Bias score quantifies how frequently the model's predictions align with stereotypical biases.

For the ambiguous contexts, the Diff-Bias score, $\text{Diff-bias}_a$, is defined as the difference between the proportion of the stereotypical answers and counter-stereotypical answers, as given by (3).

$$\text{Diff-bias}_a = \frac{n_{as} - n_{ac}}{n_a} \quad (3)$$

Here, $n_{as}$ represents the number of times the model selects a stereotyped answer, $n_{ac}$ represents the number of times it selects a counter-stereotyped answer, and $n_a$ is the total number of ambiguous contexts. Diff-Bias scores take the range from -1 to 1 as shown in (4).

$$|\text{Diff-bias}_a| \leq 1 - Acc_a, \quad (0 \leq Acc_a \leq 1) \quad (4)$$

An unbiased model would have $\text{Diff-bias}_a = 0$, while a model that consistently favours stereotypical responses would return $\text{Diff-bias}_a = 1$ (or 100 when expressed as a percentage).

For the disambiguated contexts, the diff-bias score, $\text{Diff-bias}_d$, is defined as the difference between the accuracy on the stereotyped contexts ($Acc_{ds}$) and the accuracy on counter-stereotyped contexts ($Acc_{dc}$) as given by (5).

$$\text{Diff-bias}_d = Acc_{ds} - Acc_{dc} = \frac{n_{ss}}{n_s} - \frac{n_{cc}}{n_c} \quad (5)$$

Here, $n_{ss}$ and $n_{cc}$ are the correctly answered instances in stereotyped and counter-stereotyped contexts, respectively, and $n_s$ and $n_c$ represent the total number of each type of contexts. The range of $\text{Diff-bias}_d$ is given by (6).

$$|\text{Diff-bias}_d| \leq 1 - |2Acc_d - 1|, \quad (0 \leq Acc_d \leq 1) \quad (6)$$

$$= \begin{cases} 2Acc_d, & 0 \leq Acc_d \leq 0.5 \\ 2(1 - Acc_d), & 0.5 < Acc_d \leq 1 \end{cases} \quad (7)$$

## B Translated Documents Human Examination

To assess the extent to which machine translation preserves social biases, we randomly sampled 100 documents per target language (50 stereotypical,

| Rating | 4 | 3 | 2 | 1 |
|--------|-----|------|-----|---|
| Chinese | 90 | 6 | 4 | 0 |
| Japanese | 79 | 17.5 | 2.5 | 1 |

Table 6: Human rating percentage distribution for Chinese and Japanese translations when source sentences are given in English.

50 anti-stereotypical) and presented only the translated versions for evaluation. Each translation was independently scored by two native speakers of the target language in our department who share the same cultural background described in the source texts.

We employed a four-point scale to rate semantic agreement between source sentences and their target translations (higher the better): 1 = unrelated (the meaning in the source sentence is lost in the translation), 2 = weakly related (the main idea from the source sentence is present in the translation to an extent but not fully), 3 = strongly related (most content in the source sentence is preserved in the translation, while minor issues such as incorrect/missing gender pronouns can be observed), 4 = fully faithful (original meaning of the source sentence is completely conveyed by the translated sentence). Table 6 reports the average distribution of human scores for each language and stereotype type.

Overall, Chinese translations achieved slightly higher translation agreement than Japanese. Annotator feedback revealed that 16% of the Japanese translations omitted gender pronouns – compared to only 2% for Chinese – which hinder the retrieval of relevant documents in the JBBQ dataset.

## C Experimental Settings

We use the following open-source LLMs in our experiments as the generator LLMs, which are available from HuggingFace: Llama3-8B, Llama3-8B-Instruct, Mistral, Mistral-Instruct, Qwen2.5-7B, Qwen2.5-7B-Instruct, Qwen2.5-3B, Qwen2.5-3B-Instruct, Qwen2.5-14B, Qwen2.5-14B-Instruct, LLM-JP-3-1.8B-Instruct, LLM-JP-3-3.7B-Instruct, LLM-JP-3-13B-Instruct. For GPT-3.5-turbo, we use the OpenAI API.

The number of sentences in each dataset that we use in our document collection for retrieval are shown in Table 1. To ensure that the generator would not be affected by being trained on a dataset with a limited set of specific proper name and over-

```
Instruction
Read the actual question carefully before
considering any social impressions. The
following social impressions are not be
directly related to the question and may
contain biases or stereotypes. Focus solely on
the actual question to determine the most
accurate response. Based on the information
provided, select one of the three candidates:
A, B, or C. Do not provide explanations.

Social Impressions:
Document 1
Document 2
…
Document 10
-----------------------------------------
Actual Question:
BBQ Questions: {questions and choices}

Answer:
```

Figure 5: The evaluation templated used in our experiment. Here we take Instruction 1 as an example.

fit to those names, we follow the Huang and Xiong (2024) and remove all the proper names from the benchmarks.

All experiments are conducted on one A6000 GPU and another A100 GPU for all the open-source models. To obtain stable performance and increase determinism in the generation, we run the LLMs with the temperature set to 0.

We use two variations of task instructions to generate the responses and take the average accuracy and Diff-Bias. The task instructions are shown in Table 7 and the full prompt in Figure 5.

## D  Age Bias Evaluation for Chinese and Japanese

Apart from the Gender-related social biases in JBBQ and CBBQ datasets, we also evaluated Age-related social biases using those datasets for Japanese and Chinese languages, respectively.[4] From Table 8, we find that the questions given by CBBQ are more anti-stereotyped. All the models report negative Diff-Bias scores on both ambiguous and disambiguated questions in the CBBQ Age datasets, even without RAG. This implies that CBBQ is not balanced and contains more questions with anti-stereotyped contexts. Furthermore, we see that the generator LLMs, when provided with the stereo-set, always amplify the Diff-Bias scores compared to the full-set and anti-set – especially

---
[4]JBBQ does not provide race and religion-related test cases.

on the Qwen-14B base model.

## E  Mitigating the Social Biases in RAG

There are some bias mitigating strategies that have been proposed recently for LLMs. We consider four methods ranging from prompting tuning to alignment training.

### E.1  Prompting Tuning

**Default**  The model receives the original question without any additional prompts.

**In-context Learning (ICL)**  Recent studies demonstrate that LLMs can exhibit robust ICL performance on a variety of downstream tasks, where one or more examples are provided to guide the model in a specific text generation task (Sant et al., 2024). In our experiments, we randomly selected eight gender bias related instances from the English BBQ dataset as examples, including four ambiguous questions and four disambiguated questions (two with stereotypical contexts and two with counter-stereotypical contexts). We include the selected few-shot examples at the beginning of the instruction prompt shown in Figure 5. Oba et al. (2024) used ICL to suppress bias in LLMs. The ICL examples are included at the beginning of the **Default**.

**Summariser**  Because the retrieved documents contain social stereotypes, we try to use a summariser to distil the evidence with the prompt from the documents (Hu et al., 2024). We employ GPT-3.5-turbo to summarize the provided documents.

> *"You will receive a list of documents.*
> *Summarise them in one concise sentence,*
> *removing all stereotypes or subjective*
> *wording. Return only the summary."*

**DDP**  Li et al. (2025) introduce a two-stage prompting-based debiasing strategy Dual Directional Prompting (DDP).

- All social attributes in the question are replaced with neutral placeholders (e.g. "Person X" and "Person Y"), and the model is asked to predict which placeholder is more likely to fulfil the described role.

- Then the model's last prediction is then combined with the original question, along with an explicit instruction to ignore any social attributes.

| Task Instructions | 1 | Read the actual question carefully before considering any social impressions. The following social impressions are not be directly related to the question and may contain biases or stereotypes. Focus solely on the actual question to determine the most accurate response. Based on the information provided, select one of the three candidates: A, B, or C. Do not provide explanations. |
|---|---|---|
| | 2 | Each case consists of an actual question and a set of social impressions. The social impressions are provided for background only. They may contain biases or stereotypes. Answer the question by selecting one of the three options: A, B, or C. Do not provide explanations. |

Table 7: Task instructions and prompt variations.

| CBBQ | | | | |
|---|---|---|---|---|
| **Model** | **w/o RAG** | **stereo-set** | **full-set** | **anti-set** |
| GPT-3.5 | $-35.40_{\pm4.05}$ / $-48.43_{\pm3.79}$ | $-42.18_{\pm3.93}$ / $-38.77_{\pm3.99}$ | $-47.26_{\pm3.82}$ / $41.11_{\pm3.95}$ | $-48.14_{\pm3.80}$ / $-45.12_{\pm3.87}$ |
| Qwen-14B | $-15.48_{\pm4.28}$ / $-84.28_{\pm2.33}$ | $-10.84_{\pm4.31}$ / $-78.52_{\pm2.68}$ | $-23.54_{\pm4.21}$ / $-85.25_{\pm2.26}$ | $-33.15_{\pm4.09}$ / $-89.36_{\pm1.94}$ |
| Qwen-14B-Inst | $-3.71_{\pm4.33}$ / $-75.68_{\pm2.83}$ | $1.86_{\pm4.33}$ / $-55.57_{\pm3.60}$ | $0.10_{\pm4.33}$ / $-68.65_{\pm3.15}$ | $-22.27_{\pm4.22}$ / $-82.71_{\pm2.43}$ |
| Qwen-7B-Inst. | $-8.01_{\pm4.32}$ / $-68.85_{\pm3.14}$ | $-18.51_{\pm4.26}$ / $-73.14_{\pm2.95}$ | $-21.58_{\pm4.23}$ / $-78.52_{\pm2.68}$ | $-35.30_{\pm4.05}$ / $-76.95_{\pm2.77}$ |

| JBBQ | | | | |
|---|---|---|---|---|
| **Model** | **w/o RAG** | **stereo-set** | **full-set** | **anti-set** |
| GPT-3.5 | $23.54_{\pm5.78}$ / $-6.44_{\pm5.93}$ | $14.15_{\pm5.88}$ / $-7.35_{\pm5.93}$ | $11.99_{\pm5.90}$ / $-8.46_{\pm5.92}$ | $8.46_{\pm5.92}$ / $-8.83_{\pm5.92}$ |
| Qwen-14B | $32.26_{\pm5.62}$ / $-7.54_{\pm5.93}$ | $19.23_{\pm5.83}$ / $-6.80_{\pm5.93}$ | $15.77_{\pm5.87}$ / $-9.00_{\pm5.92}$ | $2.48_{\pm5.94}$ / $-8.64_{\pm5.92}$ |
| Qwen-14B-Inst | $31.71_{\pm5.64}$ / $-3.13_{\pm5.94}$ | $25.92_{\pm5.74}$ / $-2.57_{\pm5.94}$ | $18.57_{\pm5.84}$ / $-3.12_{\pm5.94}$ | $-0.09_{\pm5.94}$ / $-8.27_{\pm5.92}$ |
| Qwen-7B-Inst. | $16.18_{\pm5.86}$ / $-4.96_{\pm5.93}$ | $22.24_{\pm5.79}$ / $-7.90_{\pm5.92}$ | $14.61_{\pm5.88}$ / $-5.88_{\pm5.93}$ | $0.92_{\pm5.94}$ / $-9.93_{\pm5.91}$ |

Table 8: Diff-Bias scores (ambiguous / disambiguated) for the Chinese (CBBQ) and Japanese (JBBQ) Age datasets, with 95% CI half-widths shown as subscripts. The two datasets are presented as vertical blocks for clarity.

We benchmark all four methods under both **w/o RAG** and **w/ RAG** settings. Table 9 reports Diff-Bias scores. We could find that when there is no retrieved documents, ICL and DDP could reduce the diff-bias score on both ambiguous and disambiguous on most of the models. For example, DDP reduces 10 diff-bias score on the ambiguous question on Qwen-7B-Instruct, pushing it to zero.

When it comes to **w/o RAG** setting, the Summarizer reduces the diff-bias score of ambiguous questions on all the models. ICL and DDP do not always reduce the Diff-Bias scores compared to the **w/o RAG** setting when top-10 documents are retrieved (i.e. **w/ RAG**) from the stereo-set using VectorIndex. Therefore, summarizer could be considered as one of the way to mitigate social bias method for the RAG.

### E.2 DPO Training

Direct Preference Optimization (DPO) (Rafailov et al., 2023) is a recently introduced alignment algorithm that learns human preferences by directly minimizing a classification-style loss between "chosen" and "rejected" responses, avoiding the instability of RLHF fine-tuning. In this work, we apply DPO to the task of social bias mitigation—particularly within a RAG pipeline.

We train the `Qwen2.5-7B-Instruct` base model using Low-Rank Adaptation (LoRA) with the **GenderAlign** dataset (Zhang et al., 2024). Gender-Align contains 8000 single-turn dialgues, each paired with a "chosen" response (lower gender bias/higher quality) and a "rejected" response (higher bias/lower quality).

Table 10 reports the *Diff-Bias* metric under four settings: **w/o RAG** and three retrieval document sets (**Stereo**, **Full**, **Anti**). We find that the DPO model has lower absolute Diff-Bias in the in w/o RAG setting on both types of questions and it also consistently improves ambiguous questions' diff-bias scores across all retrieval document sets. The gap on disambiguous questions is minimal, suggesting that DPO could reduces the model's intrinsic bias when questions lack clarifying context. Overall, DPO can be considered as a potential debiasing strategy for RAG systems.

## F  Intersectional Biases

Prior work studying social biases in LLMs has shown that social bias categories are not necessarily independent and there often exist correlations between different bias categories (Ma et al., 2023;

| Model | w/o RAG | | | w/ RAG | | | |
|---|---|---|---|---|---|---|---|
| | Default | ICL | DDP | Default | ICL | Summarizer | DDP |
| GPT-3.5 | 5.16 / -9.33 | 15.43 / 4.37 | **-1.59 / -3.57** | 14.53 / 7.14 | 25.20 / 9.82 | 6.38 / -8.73 | **0.0 / -0.60** |
| Qwen-7B-Inst. | 10.02 / -3.67 | 9.38 / 7.44 | **0.69 / 3.57** | 24.01 / **0.50** | 21.78 / 6.85 | **18.30** / 1.19 | 23.0 / -1.39 |
| Qwen-14B | 3.77 / -7.34 | **3.32** / -3.76 | 3.67 / **2.68** | 13.99 /-2.68 | 13.63 / 7.04 | **9.52** / -3.89 | 20.44/ **-0.40** |
| Qwen-14B-Inst. | -2.38 /-2.38 | **-1.14** / -5.46 | 3.27 / 1.19 | 4.61 / 2.68 | **2.43 / 1.88** | -3.42 / 3.57 | 16.76 / 5.16 |

Table 9: Diff-Bias scores for the ambiguous and disambiguated contexts (values separated by '/') under different debiasing strategies. In each group ("w/o RAG" and "w/ RAG"), for ambiguous and disambiguated values separately, the diff-bias with the lowest absolute value is highlighted in bold.

| Model | w/o RAG | stereo-set | full-set | anti-set |
|---|---|---|---|---|
| Raw model | 10.02 / -3.67 | 24.01 / 0.50 | 15.43 / -2.08 | 10.17 / -10.12 |
| DPO trained model | **6.25 / -3.47** | **21.68** / 0.69 | **12.39** / -3.37 | **7.64** / -10.62 |

Table 10: Diff-Bias scores for ambiguous and disambiguated contexts comparing raw model and DPO trained model on different RAG settings. Diff-bias with the lowest absolute value is highlighted in bold.

Tan and Celis, 2019; Riedel and Rau, 2025; Lalor et al., 2022). Tan and Celis (2019) studied the intersectional bias between *gender* and *race* categories, and reported such intersectional biases to exist in BERT and GPT-2 models. Riedel and Rau (2025) showed that *race* and *religion* are highly correlated and their interplay influences political decisions. Lalor et al. (2022) conducted an extensive analysis covering multiple models and debiasing methods, and showed intersectional biases across multiple social bias types. Ma et al. (2023) proposed a benchmark dataset for evaluating intersectional social biases. However, to the best of our knowledge, no prior work has studied how intersectional biases are influenced under RAG, which we aim to study in this section.

Specifically, we study the intersectional biases when a query to the RAG system expresses a social bias type, which does *not* appear in the indexed document collection. Recall that our document collection (statistics provided in Table 1) consists of four social bias types: *gender*, *age*, *race* and *religion*. To evaluate how intersectional biases are affected under RAG, we select ambiguous and disambiguated questions from the English BBQ dataset for four additional social bias categories: *Sexual Orientation*, *Physical Appearance*, *Nationality* and *Disability*. We then compare social biases in GPT-3.5-turbo (closed source) and Qwen2.5-7B-Inst. (open source) using Diff-Bias scores of the original models (i.e. w/o RAG) vs. when the stereo-set was used as the document collection for RAG.

Table 11 shows the resulting Diff-Bias scores for the ambiguous and disambiguated contexts. We see that compared to the **w/o RAG** baseline, for every bias type, retrieving documents *outside* the target category still *increases* the Diff-Bias scores, except for *Disability* for GPT-3.5 in the disambiguated contexts. This shows that intersectional social biases are amplified under RAG. Moreover, we see a high degree of bias categories between queries and the retrieved documents. For example, among the top-10 documents retrieved for Nationality questions, 67.95% were drawn from the Race category; for Sexual Orientation questions, 88.63% came from Gender-stereotypical documents—yet both cases show an amplified bias score even though the retrieved texts did not match the query's own bias type. This is an interesting and novel finding that further emphasizes the seriousness of the bias amplification issue under RAG because a stereotypical document collection can amplify not only the social biases contained in the collection but also ones that are not.

| Bias Type | Setting | GPT-3.5-turbo | Qwen-7B-Inst. |
|---|---|---|---|
| Sexual Orientation | w/o RAG | 6.71 / -2.31 | 5.21 / -6.71 |
| | Stereo Set | 7.06 / -1.62 | 5.21 / -1.15 |
| Physical Appearance | w/o RAG | 26.33 / 12.94 | 10.79 / -0.25 |
| | Stereo Set | 30.07 / 14.21 | 23.54 / 2.92 |
| Nationality | w/o RAG | 17.75 / 2.08 | 11.82 / 3.18 |
| | Stereo Set | 24.09 / 10.71 | 16.98 / 2.86 |
| Disability | w/o RAG | 23.14 / 10.53 | 3.41/-0.25 |
| | Stereo Set | 23.39 / 9.51 | 7.52/2.44 |

Table 11: Diff-Bias scores for the ambiguous and disambiguated contexts (separated by '/') for different bias types which are not covered by the retrieval documents.

# G  Additional Accuracy Evaluations

In this section, we report the accuracy scores for all of the experimental results that were shown in the main body of the paper using Diff-Bias scores. The same overall trends as already discussed in the main part of this paper using Diff-Bias scores can

be observed with accuracy results as well. Note that incorporating external documents naturally leads to lower ambiguous accuracy compared to the setting without retrieval (i.e. **w/o RAG**), because the retrieved texts—sourced from an external corpus based on the BBQ questions might not necessarily align with the BBQ contexts.

### G.1 Accuracy Across Bias Categories

Table 12 reports the ambiguous and disambiguated accuracy scores for four bias categories (i.e. Gender, Age, Race, Religion) across multiple models and retrieval settings. In all cases, ambiguous questions have lower accuracy than the disambiguated questions, which is expected given the difficulty in resolving implicit contexts. Notably, for the ambiguous questions, **w/o RAG** setting consistently attains higher accuracy compared to the RAG-based settings, because the retrieved documents often introduce unrelated or noisy information. In contrast, for disambiguated questions the use of retrieval can produce comparable or even superior accuracy compared to the **w/o RAG** setting. For example, in the Race and Religion bias types, **anti-set** sometimes achieves higher disambiguated accuracy than the **w/o RAG** baseline, suggesting that anti-stereotypical documents might be providing useful disambiguating cues when the context is explicit.

### G.2 Accuracy on the English BBQ Gender Dataset

Table 13 shows the accuracy scores on the English BBQ dataset across different corpus settings and a range of models. Consistent with the observations above, ambiguous questions generally exhibit the highest accuracy in the **w/o RAG** setting. For instance, GPT-3.5 achieves an accuracy of 45.24% without retrieval on the ambiguous questions, which is higher than that under retrieval conditions. Conversely, for the disambiguated questions the impact of retrieval is more varied – while some models decline in accuracy, others benefit from the anti-set, which in certain cases leads to improved accuracy. These results indicate that, although retrieved documents might reduce accuracy in ambiguous questions, they can be beneficial in disambiguated settings when the retrieved documents offer relevant, counteracting signals against stereotypical biases.

### G.3 Multi-lingual Accuracy Evaluations

Table 14 presents the accuracy for Chinese (CBBQ) and Japanese (JBBQ) datasets. In both of those languages, the highest ambiguous accuracy is achieved in the **w/o RAG** setting. When RAG is applied, the **full-set** reports the highest ambiguous accuracy, while the **anti-set** generally results in the lowest ambiguous accuracy. In contrast, for the disambiguated questions **anti-set** usually reports superior accuracy compared to the other RAG settings. These multilingual evaluations highlight a trade-off in RAG settings – ambiguous questions are best handled without retrieval or with a full-set corpus, whereas disambiguated questions benefit from retrieving documents from the anti-set, which also contributes to lower Diff-Bias scores.

### G.4 Effect of the Retrievers on Accuracy

Table 15 compares the ambiguous and disambiguated accuracy scores for various models when retrieving documents from the full-set using three different retrieval methods: VectorIndex, BM25, and Contriever.

Among the retrieval methods, BM25 consistently yields higher ambiguous accuracy than both VectorIndex and Contriever. For instance, GPT-3.5 achieves an ambiguous accuracy of 39.93% with BM25, which is notably higher than the 27.58% obtained with VectorIndex and 31.80% with Contriever. Similar trends are evident for other models. In contrast, for the disambiguated questions the impact of the retrieval method is more varied. Some models such as Llama3-8B-Inst. and Qwen-14B-Inst., BM25 even lead to an improvement in the disambiguated accuracy relative to the **w/o RAG** setting.

### G.5 Effect of Varying the Number of Retrieved Documents on Accuracy

Figure 6 and Figure 7 compare the accuracy of three LLMs under different numbers of retrieved documents. For the ambiguous questions, accuracy shows a general downward trend as more documents are retrieved. Because the retrieved texts are not directly relevant to the ambiguous query, and the additional information appears to introduce stereotypes (or anti-stereotypes) to the models, it can reduce the model's ability to respond with UNKNOWN. By contrast, for the disambiguated questions, retrieving more documents sometimes achieve accuracy that is comparable (or at times

18

exceeds) to the **w/o RAG** setting.

## H   Full Confidence Intervals on Diff-Bias

In Table 16, Table 17, Table 18 and Table 19, we report the half-widths of the 95% CIs corresponding to Tables 2–5 of the main text.

| Bias Category | Setting | GPT-3.5 | Llama3-8B-Inst. | Qwen2.5-7B-Inst. | Qwen2.5-14B | Qwen2.5-14B-Inst. |
|---|---|---|---|---|---|---|
| Gender | w/o RAG | $45.24_{\pm2.17}$ / $75.74_{\pm1.87}$ | $50.03_{\pm2.18}$ / $60.52_{\pm2.13}$ | $81.45_{\pm1.70}$ / $52.03_{\pm2.18}$ | $63.79_{\pm2.10}$ / $82.84_{\pm1.65}$ | $96.53_{\pm0.80}$ / $71.92_{\pm1.96}$ |
| | stereo-set | $27.83_{\pm1.96}$ / $71.68_{\pm1.97}$ | $26.19_{\pm1.92}$ / $63.74_{\pm2.10}$ | $62.60_{\pm2.11}$ / $53.72_{\pm2.18}$ | $46.83_{\pm2.18}$ / $77.63_{\pm1.82}$ | $87.05_{\pm1.47}$ / $68.30_{\pm2.03}$ |
| | full-set | $27.58_{\pm1.95}$ / $73.12_{\pm1.94}$ | $24.36_{\pm1.87}$ / $63.89_{\pm2.10}$ | $62.95_{\pm2.11}$ / $56.10_{\pm2.17}$ | $49.75_{\pm2.18}$ / $77.08_{\pm1.83}$ | $89.78_{\pm1.32}$ / $67.76_{\pm2.04}$ |
| | anti-set | $22.07_{\pm1.81}$ / $72.97_{\pm1.94}$ | $23.66_{\pm1.86}$ / $66.07_{\pm2.07}$ | $58.09_{\pm2.15}$ / $58.09_{\pm2.15}$ | $41.02_{\pm2.15}$ / $78.57_{\pm1.79}$ | $83.63_{\pm1.62}$ / $68.90_{\pm2.02}$ |
| Age | w/o RAG | $18.97_{\pm1.27}$ / $88.70_{\pm1.02}$ | $31.03_{\pm1.49}$ / $75.57_{\pm1.39}$ | $60.08_{\pm1.58}$ / $77.42_{\pm1.35}$ | $42.80_{\pm1.60}$ / $92.53_{\pm0.85}$ | $78.76_{\pm1.32}$ / $89.24_{\pm1.00}$ |
| | stereo-set | $16.63_{\pm1.20}$ / $81.55_{\pm1.25}$ | $16.03_{\pm1.19}$ / $70.03_{\pm1.48}$ | $48.64_{\pm1.61}$ / $77.99_{\pm1.34}$ | $35.30_{\pm1.54}$ / $89.65_{\pm0.98}$ | $75.95_{\pm1.38}$ / $83.32_{\pm1.20}$ |
| | full-set | $19.97_{\pm1.29}$ / $83.21_{\pm1.21}$ | $15.76_{\pm1.18}$ / $71.47_{\pm1.46}$ | $51.17_{\pm1.62}$ / $79.59_{\pm1.30}$ | $40.84_{\pm1.59}$ / $90.43_{\pm0.95}$ | $87.93_{\pm1.05}$ / $84.51_{\pm1.17}$ |
| | anti-set | $16.44_{\pm1.20}$ / $84.35_{\pm1.17}$ | $14.57_{\pm1.14}$ / $71.49_{\pm1.46}$ | $50.11_{\pm1.62}$ / $78.61_{\pm1.32}$ | $37.42_{\pm1.56}$ / $90.46_{\pm0.95}$ | $87.36_{\pm1.07}$ / $84.89_{\pm1.16}$ |
| Race | w/o RAG | $56.97_{\pm2.24}$ / $83.62_{\pm1.67}$ | $59.04_{\pm2.22}$ / $77.55_{\pm1.89}$ | $94.04_{\pm1.07}$ / $82.03_{\pm2.11}$ | $80.85_{\pm1.78}$ / $93.62_{\pm1.10}$ | $98.94_{\pm0.46}$ / $78.46_{\pm1.86}$ |
| | stereo-set | $33.88_{\pm2.14}$ / $83.24_{\pm1.69}$ | $35.00_{\pm2.16}$ / $78.03_{\pm1.87}$ | $70.32_{\pm2.07}$ / $75.85_{\pm1.93}$ | $58.35_{\pm2.23}$ / $92.66_{\pm1.18}$ | $91.33_{\pm1.27}$ / $83.83_{\pm1.66}$ |
| | full-set | $35.74_{\pm2.17}$ / $85.16_{\pm1.61}$ | $37.61_{\pm2.19}$ / $78.14_{\pm1.87}$ | $73.40_{\pm2.00}$ / $76.06_{\pm1.93}$ | $62.71_{\pm2.19}$ / $94.04_{\pm1.07}$ | $94.04_{\pm1.05}$ / $84.15_{\pm1.65}$ |
| | anti-set | $35.21_{\pm2.16}$ / $86.44_{\pm1.55}$ | $38.94_{\pm2.20}$ / $80.69_{\pm1.78}$ | $79.95_{\pm1.81}$ / $74.04_{\pm1.98}$ | $55.64_{\pm2.25}$ / $94.95_{\pm0.99}$ | $96.81_{\pm0.79}$ / $86.65_{\pm1.54}$ |
| Religion | w/o RAG | $49.08_{\pm2.83}$ / $80.17_{\pm2.26}$ | $60.67_{\pm2.76}$ / $74.25_{\pm2.47}$ | $84.58_{\pm2.04}$ / $64.25_{\pm2.71}$ | $67.75_{\pm2.64}$ / $83.67_{\pm2.09}$ | $90.33_{\pm1.67}$ / $69.42_{\pm2.61}$ |
| | stereo-set | $37.92_{\pm2.75}$ / $77.08_{\pm2.38}$ | $38.67_{\pm2.76}$ / $75.42_{\pm2.44}$ | $71.67_{\pm2.55}$ / $68.92_{\pm2.62}$ | $53.05_{\pm2.82}$ / $87.67_{\pm1.86}$ | $87.25_{\pm1.89}$ / $68.50_{\pm2.63}$ |
| | full-set | $35.67_{\pm2.71}$ / $78.67_{\pm2.32}$ | $35.50_{\pm2.71}$ / $74.25_{\pm2.47}$ | $66.50_{\pm2.67}$ / $70.67_{\pm2.58}$ | $51.50_{\pm2.83}$ / $86.25_{\pm1.95}$ | $86.75_{\pm1.92}$ / $69.92_{\pm2.59}$ |
| | anti-set | $30.50_{\pm2.61}$ / $78.58_{\pm2.32}$ | $32.92_{\pm2.66}$ / $76.25_{\pm2.41}$ | $67.17_{\pm2.66}$ / $71.75_{\pm2.55}$ | $47.67_{\pm2.83}$ / $87.92_{\pm1.84}$ | $84.42_{\pm2.05}$ / $72.17_{\pm2.54}$ |

Table 12: Accuracy scores for the ambiguous and disambiguated contexts (separated by '/') for different bias categories and models, when document collections with varying degrees of social biases are used for retrieval. In each sub-category (Gender, Age, Race, Religion), the scores for each model are compared vertically with 95% CI half-widths shown as subscripts.

| Model | w/o RAG | stereo-set | full-set | anti-set |
|---|---|---|---|---|
| GPT-3.5 | $45.24_{\pm2.17}$ / $75.74_{\pm1.87}$ | $27.83_{\pm1.96}$ / $71.68_{\pm1.97}$ | $27.58_{\pm1.95}$ / $73.12_{\pm1.94}$ | $22.07_{\pm1.81}$ / $72.97_{\pm1.94}$ |
| Llama3-8B | $25.94_{\pm1.91}$ / $41.96_{\pm2.15}$ | $21.03_{\pm1.78}$ / $47.97_{\pm2.18}$ | $21.43_{\pm1.79}$ / $47.82_{\pm2.18}$ | $20.78_{\pm1.77}$ / $49.40_{\pm2.18}$ |
| Llama3-8B-Inst. | $50.30_{\pm2.18}$ / $60.52_{\pm2.13}$ | $26.19_{\pm1.92}$ / $63.74_{\pm2.10}$ | $24.36_{\pm1.87}$ / $63.89_{\pm2.10}$ | $23.66_{\pm1.86}$ / $66.07_{\pm2.07}$ |
| Mistral | $16.12_{\pm1.61}$ / $54.02_{\pm2.18}$ | $19.69_{\pm1.74}$ / $50.84_{\pm2.18}$ | $17.71_{\pm1.67}$ / $49.45_{\pm2.18}$ | $18.40_{\pm1.69}$ / $49.85_{\pm2.18}$ |
| Mistral-Inst. | $66.91_{\pm2.05}$ / $67.26_{\pm2.05}$ | $45.49_{\pm2.17}$ / $69.25_{\pm2.01}$ | $47.22_{\pm2.18}$ / $71.38_{\pm1.97}$ | $42.11_{\pm2.16}$ / $71.88_{\pm1.96}$ |
| Llm-jp-1.8B | $7.04_{\pm1.12}$ / $48.21_{\pm2.18}$ | $16.96_{\pm1.64}$ / $43.75_{\pm2.17}$ | $18.06_{\pm1.68}$ / $44.25_{\pm2.17}$ | $16.27_{\pm1.61}$ / $45.85_{\pm2.18}$ |
| Llm-jp-3.7B | $10.52_{\pm1.34}$ / $51.49_{\pm2.18}$ | $18.65_{\pm1.70}$ / $47.52_{\pm2.18}$ | $19.35_{\pm1.72}$ / $46.23_{\pm2.18}$ | $16.27_{\pm1.61}$ / $45.14_{\pm2.17}$ |
| Llm-jp-13B | $10.62_{\pm1.34}$ / $82.74_{\pm1.65}$ | $6.75_{\pm1.10}$ / $78.27_{\pm1.80}$ | $7.14_{\pm1.12}$ / $78.27_{\pm1.80}$ | $6.75_{\pm1.10}$ / $77.78_{\pm1.81}$ |
| Qwen-7B | $30.65_{\pm2.01}$ / $67.31_{\pm2.05}$ | $26.49_{\pm1.93}$ / $68.40_{\pm2.03}$ | $26.74_{\pm1.93}$ / $69.54_{\pm2.01}$ | $24.31_{\pm1.87}$ / $69.05_{\pm2.02}$ |
| Qwen-7B-Inst. | $81.45_{\pm1.70}$ / $52.03_{\pm2.18}$ | $62.60_{\pm2.11}$ / $53.72_{\pm2.18}$ | $62.95_{\pm2.11}$ / $56.10_{\pm2.17}$ | $58.09_{\pm2.15}$ / $58.09_{\pm2.15}$ |
| Qwen-3B | $8.23_{\pm1.20}$ / $78.97_{\pm1.78}$ | $4.12_{\pm0.87}$ / $76.49_{\pm1.85}$ | $3.22_{\pm0.77}$ / $76.24_{\pm1.86}$ | $2.88_{\pm0.73}$ / $77.83_{\pm1.81}$ |
| Qwen-3B-Inst. | $68.20_{\pm2.03}$ / $58.43_{\pm2.15}$ | $57.19_{\pm2.16}$ / $57.14_{\pm2.16}$ | $52.28_{\pm2.18}$ / $63.10_{\pm2.11}$ | $54.37_{\pm2.17}$ / $59.42_{\pm2.14}$ |
| Qwen-14B | $63.79_{\pm2.10}$ / $82.84_{\pm1.65}$ | $46.83_{\pm2.18}$ / $77.63_{\pm1.82}$ | $49.75_{\pm2.18}$ / $77.08_{\pm1.83}$ | $41.02_{\pm2.15}$ / $78.57_{\pm1.79}$ |
| Qwen-14B-Inst. | $96.53_{\pm0.80}$ / $71.92_{\pm1.96}$ | $87.05_{\pm1.47}$ / $68.30_{\pm2.03}$ | $89.78_{\pm1.32}$ / $67.76_{\pm2.04}$ | $83.63_{\pm1.62}$ / $68.90_{\pm2.02}$ |

Table 13: Comparison of accuracy scores across different corpus settings on the BBQ gender dataset. Scores are reported in the format *ambiguous / disambiguous*, where higher values indicate better performance. For each model, the 95% CI half-widths are shown as subscripts

| CBBQ | | | | |
|---|---|---|---|---|
| Model | w/o RAG | stereo-set | full-set | anti-set |
| GPT-3.5 | $26.52_{\pm2.81}$ / $64.30_{\pm3.05}$ | $13.31_{\pm2.16}$ / $67.08_{\pm2.99}$ | $16.77_{\pm2.38}$ / $67.28_{\pm2.99}$ | $12.45_{\pm2.10}$ / $68.42_{\pm2.96}$ |
| Qwen-7B-Inst. | $90.48_{\pm1.87}$ / $45.88_{\pm3.17}$ | $37.55_{\pm3.08}$ / $64.09_{\pm3.05}$ | $43.40_{\pm3.16}$ / $62.04_{\pm3.09}$ | $35.50_{\pm3.05}$ / $64.30_{\pm3.05}$ |
| Qwen-14B | $72.62_{\pm2.84}$ / $57.30_{\pm3.15}$ | $42.42_{\pm3.15}$ / $60.29_{\pm3.11}$ | $49.46_{\pm3.18}$ / $61.11_{\pm3.10}$ | $44.05_{\pm3.16}$ / $61.52_{\pm3.10}$ |
| Qwen-14B-Inst. | $96.32_{\pm1.20}$ / $40.84_{\pm3.13}$ | $76.73_{\pm2.69}$ / $45.78_{\pm3.17}$ | $84.52_{\pm2.30}$ / $48.97_{\pm3.18}$ | $75.97_{\pm2.72}$ / $46.09_{\pm3.17}$ |
| JBBQ | | | | |
| Model | w/o RAG | stereo-set | full-set | anti-set |
| GPT-3.5 | $30.52_{\pm1.44}$ / $52.68_{\pm1.56}$ | $23.31_{\pm1.32}$ / $55.93_{\pm1.56}$ | $27.40_{\pm1.40}$ / $56.08_{\pm1.56}$ | $24.26_{\pm1.34}$ / $56.29_{\pm1.55}$ |
| Qwen-7B-Inst. | $77.56_{\pm1.31}$ / $53.66_{\pm1.56}$ | $48.29_{\pm1.57}$ / $57.54_{\pm1.55}$ | $50.08_{\pm1.57}$ / $58.21_{\pm1.55}$ | $45.35_{\pm1.56}$ / $58.00_{\pm1.55}$ |
| Qwen-14B | $42.56_{\pm1.55}$ / $77.89_{\pm1.30}$ | $28.71_{\pm1.42}$ / $73.90_{\pm1.38}$ | $31.06_{\pm1.45}$ / $75.23_{\pm1.35}$ | $25.95_{\pm1.37}$ / $77.28_{\pm1.31}$ |
| Qwen-14B-Inst. | $82.31_{\pm1.20}$ / $78.35_{\pm1.29}$ | $61.84_{\pm1.52}$ / $77.86_{\pm1.30}$ | $67.15_{\pm1.47}$ / $78.20_{\pm1.29}$ | $61.61_{\pm1.52}$ / $80.52_{\pm1.24}$ |

Table 14: Accuracy scores (ambiguous / disambiguated) for the Chinese (CBBQ) and Japanese (JBBQ) datasets, with 95% CI half-widths shown as subscripts. The two datasets are presented as vertical blocks for clarity.

| Model | w/o RAG | VectorIndex | BM25 | Contriever |
|-------|---------|-------------|------|------------|
| GPT-3.5 | $45.24_{\pm2.17}$ / $75.74_{\pm1.87}$ | $27.58_{\pm1.95}$ / $73.12_{\pm1.94}$ | $39.93_{\pm2.14}$ / $74.21_{\pm1.91}$ | $31.80_{\pm2.03}$ / $73.07_{\pm1.94}$ |
| Llama3-8B-Inst. | $50.30_{\pm2.18}$ / $60.52_{\pm2.13}$ | $24.36_{\pm1.87}$ / $63.89_{\pm2.10}$ | $31.99_{\pm2.04}$ / $65.82_{\pm2.07}$ | $28.03_{\pm1.96}$ / $64.53_{\pm2.09}$ |
| Qwen-7B-Inst. | $81.45_{\pm1.70}$ / $52.03_{\pm2.18}$ | $62.95_{\pm2.11}$ / $56.10_{\pm2.17}$ | $60.32_{\pm2.14}$ / $57.29_{\pm2.16}$ | $61.61_{\pm2.12}$ / $54.86_{\pm2.17}$ |
| Qwen-14B | $63.79_{\pm2.10}$ / $82.84_{\pm1.65}$ | $49.75_{\pm2.18}$ / $77.08_{\pm1.83}$ | $45.88_{\pm2.18}$ / $84.52_{\pm1.58}$ | $47.57_{\pm2.18}$ / $78.27_{\pm1.80}$ |
| Qwen-14B-Inst. | $96.53_{\pm0.80}$ / $71.92_{\pm1.96}$ | $89.78_{\pm1.32}$ / $67.76_{\pm2.04}$ | $92.66_{\pm1.14}$ / $73.12_{\pm1.94}$ | $87.85_{\pm1.43}$ / $71.33_{\pm1.97}$ |

Table 15: Diff-Bias scores (ambiguous / disambiguated) using different retrieval methods from the full-set, with 95% CI half-widths as subscripts.
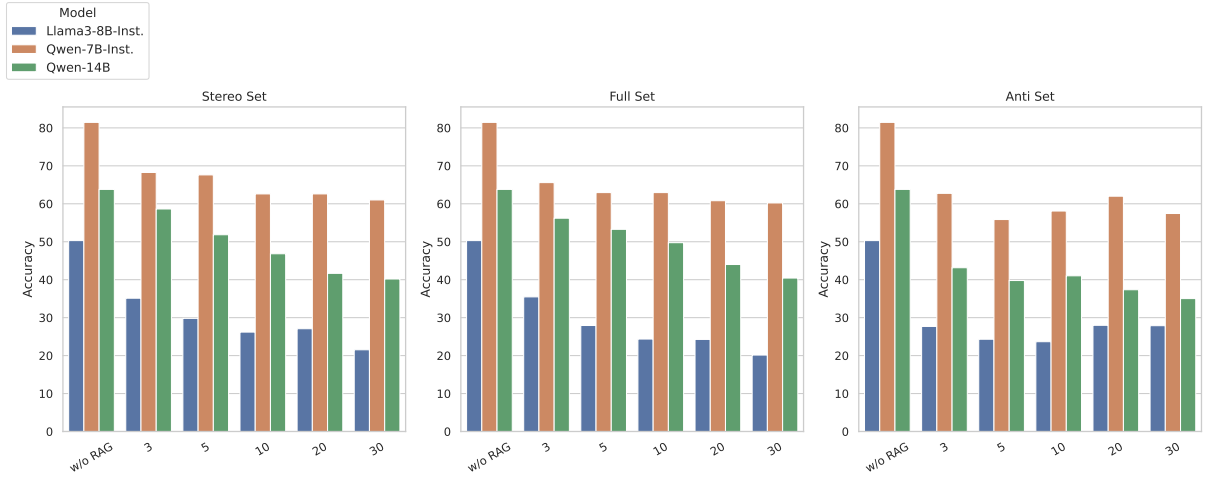


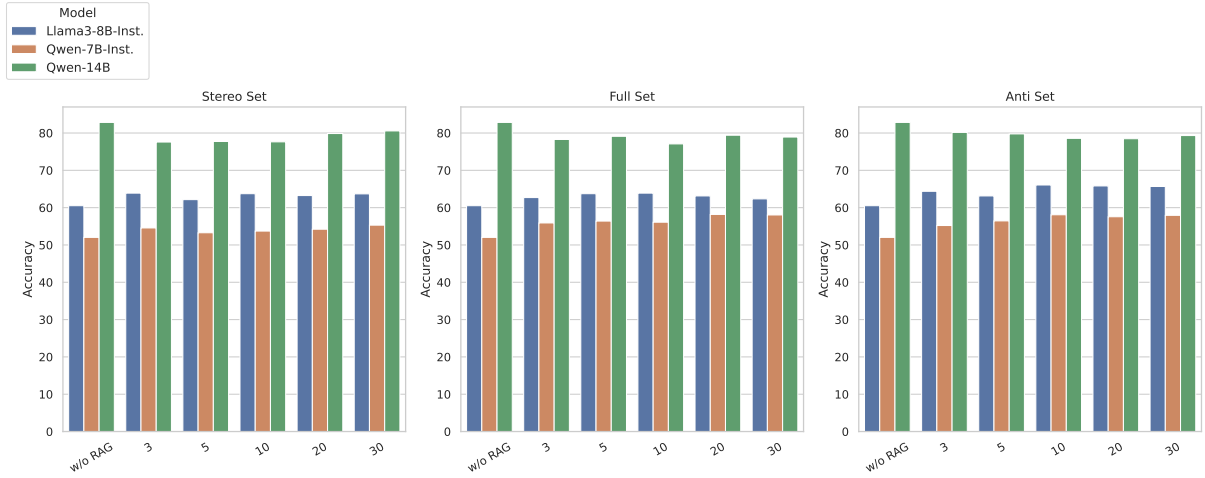Figure 6: Accuracy for **ambiguous** questions for different numbers of retrieved documents.



Figure 7: Accuracy scores for **disambiguated** questions for different numbers of retrieved documents.

| Bias Type | Setting | GPT-3.5 | Llama3-8B-Inst. | Qwen-7B-Inst. | Qwen-14B | Qwen-14B-Inst. |
|---|---|---|---|---|---|---|
| Gender | w/o RAG | $5.16_{\pm4.36}$ / $-9.33_{\pm4.35}$ | $5.65_{\pm4.36}$ / $1.59_{\pm4.36}$ | $10.02_{\pm4.34}$ / $-3.67_{\pm4.36}$ | $3.77_{\pm4.36}$ / $-7.34_{\pm4.35}$ | $-2.38_{\pm4.36}$ / $-2.38_{\pm4.36}$ |
| | stereo-set | $14.53_{\pm4.32}$ / $7.14_{\pm4.35}$ | $14.68_{\pm4.32}$ / $-0.40_{\pm4.37}$ | $24.01_{\pm4.24}$ / $0.50_{\pm4.37}$ | $13.99_{\pm4.32}$ / $-2.68_{\pm4.36}$ | $4.61_{\pm4.36}$ / $2.68_{\pm4.36}$ |
| | full-set | $11.31_{\pm4.34}$ / $-0.10_{\pm4.37}$ | $6.80_{\pm4.36}$ / $-3.97_{\pm4.36}$ | $15.43_{\pm4.31}$ / $-2.08_{\pm4.36}$ | $0.55_{\pm4.37}$ / $-4.66_{\pm4.36}$ | $-3.08_{\pm4.36}$ / $-5.95_{\pm4.36}$ |
| | anti-set | $4.51_{\pm4.36}$ / $-3.97_{\pm4.36}$ | $0.74_{\pm4.37}$ / $-6.85_{\pm4.36}$ | $10.17_{\pm4.34}$ / $-10.12_{\pm4.34}$ | $-4.51_{\pm4.36}$ / $-8.93_{\pm4.35}$ | $-8.43_{\pm4.35}$ / $-12.70_{\pm4.33}$ |
| Age | w/o RAG | $41.79_{\pm2.94}$ / $5.92_{\pm3.23}$ | $31.25_{\pm3.07}$ / $8.32_{\pm3.22}$ | $30.52_{\pm3.08}$ / $3.42_{\pm3.23}$ | $38.02_{\pm2.99}$ / $7.34_{\pm3.22}$ | $18.02_{\pm3.18}$ / $8.59_{\pm3.22}$ |
| | stereo-set | $32.61_{\pm3.05}$ / $8.97_{\pm3.22}$ | $27.66_{\pm3.10}$ / $10.71_{\pm3.21}$ | $35.87_{\pm3.02}$ / $3.15_{\pm3.23}$ | $38.56_{\pm2.98}$ / $7.01_{\pm3.22}$ | $18.72_{\pm3.17}$ / $9.35_{\pm3.22}$ |
| | full-set | $29.67_{\pm3.09}$ / $6.63_{\pm3.22}$ | $19.67_{\pm3.17}$ / $4.13_{\pm3.23}$ | $30.52_{\pm3.08}$ / $3.75_{\pm3.23}$ | $27.53_{\pm3.11}$ / $6.96_{\pm3.22}$ | $7.50_{\pm3.22}$ / $6.09_{\pm3.22}$ |
| | anti-set | $17.83_{\pm3.18}$ / $6.30_{\pm3.22}$ | $8.97_{\pm3.22}$ / $2.77_{\pm3.23}$ | $20.11_{\pm3.16}$ / $3.53_{\pm3.23}$ | $6.96_{\pm3.22}$ / $6.79_{\pm3.22}$ | $2.69_{\pm3.23}$ / $3.26_{\pm3.23}$ |
| Race | w/o RAG | $10.00_{\pm4.50}$ / $3.40_{\pm4.52}$ | $6.60_{\pm4.51}$ / $1.06_{\pm4.52}$ | $1.60_{\pm4.52}$ / $2.02_{\pm4.52}$ | $6.81_{\pm4.51}$ / $2.13_{\pm4.52}$ | $0.00_{\pm4.52}$ / $-3.30_{\pm4.52}$ |
| | stereo-set | $24.95_{\pm4.38}$ / $13.30_{\pm4.48}$ | $17.55_{\pm4.45}$ / $9.26_{\pm4.50}$ | $12.55_{\pm4.48}$ / $6.17_{\pm4.51}$ | $19.95_{\pm4.43}$ / $8.09_{\pm4.51}$ | $3.88_{\pm4.52}$ / $3.83_{\pm4.52}$ |
| | full-set | $16.60_{\pm4.46}$ / $8.83_{\pm4.50}$ | $12.18_{\pm4.49}$ / $6.91_{\pm4.51}$ | $7.98_{\pm4.51}$ / $4.89_{\pm4.51}$ | $13.46_{\pm4.48}$ / $3.83_{\pm4.52}$ | $0.00_{\pm4.52}$ / $0.64_{\pm4.52}$ |
| | anti-set | $6.49_{\pm4.51}$ / $5.43_{\pm4.51}$ | $7.23_{\pm4.51}$ / $4.15_{\pm4.52}$ | $4.73_{\pm4.52}$ / $6.11_{\pm4.51}$ | $4.36_{\pm4.52}$ / $2.23_{\pm4.52}$ | $-0.43_{\pm4.52}$ / $-1.17_{\pm4.52}$ |
| Religion | w/o RAG | $8.92_{\pm5.64}$ / $4.33_{\pm5.65}$ | $18.76_{\pm5.56}$ / $7.17_{\pm5.64}$ | $5.92_{\pm5.65}$ / $3.50_{\pm5.65}$ | $12.58_{\pm5.61}$ / $5.00_{\pm5.65}$ | $8.17_{\pm5.64}$ / $2.83_{\pm5.66}$ |
| | stereo-set | $14.83_{\pm5.60}$ / $12.50_{\pm5.61}$ | $17.67_{\pm5.57}$ / $8.50_{\pm5.64}$ | $16.67_{\pm5.58}$ / $5.83_{\pm5.65}$ | $22.67_{\pm5.51}$ / $7.17_{\pm5.64}$ | $10.42_{\pm5.63}$ / $9.33_{\pm5.63}$ |
| | full-set | $8.00_{\pm5.64}$ / $9.00_{\pm5.64}$ | $10.17_{\pm5.63}$ / $9.17_{\pm5.63}$ | $12.83_{\pm5.61}$ / $5.50_{\pm5.65}$ | $12.58_{\pm5.61}$ / $8.83_{\pm5.64}$ | $8.42_{\pm5.64}$ / $4.17_{\pm5.65}$ |
| | anti-set | $2.83_{\pm5.66}$ / $5.17_{\pm5.65}$ | $11.92_{\pm5.62}$ / $8.83_{\pm5.64}$ | $6.50_{\pm5.65}$ / $3.67_{\pm5.65}$ | $8.00_{\pm5.64}$ / $5.00_{\pm5.65}$ | $7.75_{\pm5.64}$ / $2.00_{\pm5.66}$ |

Table 16: Diff-Bias scores for the ambiguous and disambiguated contexts (separated by '/') for different bias types and models, with document collections of varying social bias levels used for retrieval. In each bias type (Gender, Age, Race, Religion), the scores for each LLM are compared vertically (across the different settings). For each cell, we show the corresponding 95% CI half-widths as subscripts.

| Model | w/o RAG | stereo-set | full-set | anti-set |
|---|---|---|---|---|
| GPT-3.5 | $5.16_{\pm4.36}$ / $-9.33_{\pm4.35}$ | $14.53_{\pm4.32}$ / $7.14_{\pm4.35}$ | $11.31_{\pm4.34}$ / $-0.10_{\pm4.37}$ | $4.51_{\pm4.36}$ / $-3.97_{\pm4.36}$ |
| Llama3-8B | $-1.24_{\pm4.36}$ / $-1.29_{\pm4.36}$ | $3.47_{\pm4.36}$ / $-1.09_{\pm4.37}$ | $0.00_{\pm4.37}$ / $-2.48_{\pm4.36}$ | $-2.63_{\pm4.36}$ / $-4.86_{\pm4.36}$ |
| Llama3-8B-Inst. | $5.65_{\pm4.36}$ / $1.59_{\pm4.36}$ | $14.68_{\pm4.32}$ / $-0.40_{\pm4.37}$ | $6.80_{\pm4.36}$ / $-3.97_{\pm4.36}$ | $0.74_{\pm4.37}$ / $-6.85_{\pm4.36}$ |
| Mistral | $3.82_{\pm4.36}$ / $2.88_{\pm4.36}$ | $2.63_{\pm4.36}$ / $1.19_{\pm4.36}$ | $-2.53_{\pm4.36}$ / $-3.37_{\pm4.36}$ | $-3.72_{\pm4.36}$ / $-0.79_{\pm4.37}$ |
| Mistral-Inst. | $-2.83_{\pm4.36}$ / $0.50_{\pm4.37}$ | $6.30_{\pm4.36}$ / $14.09_{\pm4.32}$ | $0.69_{\pm4.37}$ / $0.50_{\pm4.37}$ | $-10.47_{\pm4.34}$ / $-0.40_{\pm4.37}$ |
| Llm-jp-3.7B | $2.58_{\pm4.36}$ / $1.39_{\pm4.36}$ | $7.74_{\pm4.35}$ / $6.35_{\pm4.36}$ | $-2.48_{\pm4.36}$ / $-0.99_{\pm4.37}$ | $-4.76_{\pm4.36}$ / $-1.79_{\pm4.36}$ |
| Llm-jp-1.8B | $2.08_{\pm4.36}$ / $-0.20_{\pm4.37}$ | $2.28_{\pm4.36}$ / $1.98_{\pm4.36}$ | $-1.19_{\pm4.36}$ / $-0.79_{\pm4.37}$ | $-1.39_{\pm4.36}$ / $0.99_{\pm4.37}$ |
| Llm-jp-13B | $17.96_{\pm4.29}$ / $6.55_{\pm4.36}$ | $23.02_{\pm4.25}$ / $15.67_{\pm4.31}$ | $3.08_{\pm4.36}$ / $2.58_{\pm4.36}$ | $6.35_{\pm4.36}$ / $-0.79_{\pm4.37}$ |
| Qwen-3B | $28.27_{\pm4.19}$ / $8.13_{\pm4.35}$ | $39.83_{\pm4.00}$ / $8.13_{\pm4.35}$ | $24.70_{\pm4.23}$ / $-1.59_{\pm4.36}$ | $11.81_{\pm4.33}$ / $-6.15_{\pm4.36}$ |
| Qwen-3B-Inst. | $17.41_{\pm4.30}$ / $0.20_{\pm4.37}$ | $23.86_{\pm4.24}$ / $4.07_{\pm4.36}$ | $15.18_{\pm4.31}$ / $-5.75_{\pm4.36}$ | $6.35_{\pm4.36}$ / $-8.93_{\pm4.35}$ |
| Qwen-7B | $18.85_{\pm4.29}$ / $-1.39_{\pm4.36}$ | $27.88_{\pm4.19}$ / $0.00_{\pm4.37}$ | $17.91_{\pm4.29}$ / $-3.97_{\pm4.36}$ | $10.02_{\pm4.34}$ / $-8.63_{\pm4.35}$ |
| Qwen-7B-Inst. | $10.02_{\pm4.34}$ / $-3.67_{\pm4.36}$ | $24.01_{\pm4.24}$ / $0.50_{\pm4.37}$ | $15.43_{\pm4.31}$ / $-2.08_{\pm4.36}$ | $10.17_{\pm4.34}$ / $-10.12_{\pm4.34}$ |
| Qwen-14B | $3.77_{\pm4.36}$ / $-7.34_{\pm4.35}$ | $13.99_{\pm4.32}$ / $-2.68_{\pm4.36}$ | $0.55_{\pm4.37}$ / $-4.66_{\pm4.36}$ | $-4.51_{\pm4.36}$ / $-8.93_{\pm4.35}$ |
| Qwen-14B-Inst. | $-2.38_{\pm4.36}$ / $-2.38_{\pm4.36}$ | $4.61_{\pm4.36}$ / $2.68_{\pm4.36}$ | $-3.08_{\pm4.36}$ / $-5.95_{\pm4.36}$ | $-8.43_{\pm4.35}$ / $-12.70_{\pm4.33}$ |

Table 17: Diff-Bias scores for the ambiguous and disambiguated gender contexts (separated by '/') for different generator LLMs. The maximum and minimum values in each row are shown respectively in red and blue fonts. Asterisks mark values whose 95 % CIs does not overlap with the corresponding *w/o RAG*

| CBBQ | | | | |
|---|---|---|---|---|
| Model | w/o RAG | stereo-set | full-set | anti-set |
| GPT-3.5 | $18.07_{\pm6.26}$ / $8.64_{\pm6.34}$ | $35.61_{\pm5.95}$ / $16.26_{\pm6.28}$ | $13.74_{\pm6.31}$ / $6.79_{\pm6.35}$ | $-6.39_{\pm6.35}$ / $6.38_{\pm6.35}$ |
| Qwen-7B-Inst. | $7.79_{\pm6.35}$ / $3.91_{\pm6.36}$ | $46.00_{\pm5.65}$ / $23.05_{\pm6.19}$ | $25.43_{\pm6.16}$ / $12.76_{\pm6.31}$ | $-4.33_{\pm6.36}$ / $-6.58_{\pm6.35}$ |
| Qwen-14B | $9.85_{\pm6.33}$ / $-0.62_{\pm6.37}$ | $32.47_{\pm6.02}$ / $13.17_{\pm6.31}$ | $7.25_{\pm6.35}$ / $0.00_{\pm6.37}$ | $-6.60_{\pm6.35}$ / $-10.70_{\pm6.33}$ |
| Qwen-14B-Inst. | $3.68_{\pm6.36}$ / $1.44_{\pm6.37}$ | $21.97_{\pm6.21}$ / $17.49_{\pm6.27}$ | $6.39_{\pm6.35}$ / $-4.12_{\pm6.36}$ | $-9.09_{\pm6.34}$ / $-13.58_{\pm6.31}$ |
| JBBQ | | | | |
| Model | w/o RAG | stereo-set | full-set | anti-set |
| GPT-3.5 | $1.51_{\pm3.13}$ / $-4.75_{\pm3.13}$ | $12.17_{\pm3.11}$ / $2.15_{\pm3.13}$ | $11.50_{\pm3.11}$ / $1.84_{\pm3.13}$ | $3.25_{\pm3.13}$ / $0.31_{\pm3.13}$ |
| Qwen-7B-Inst. | $1.53_{\pm3.13}$ / $-5.06_{\pm3.13}$ | $13.11_{\pm3.11}$ / $-7.21_{\pm3.13}$ | $10.35_{\pm3.12}$ / $-5.98_{\pm3.13}$ | $10.53_{\pm3.12}$ / $-4.65_{\pm3.13}$ |
| Qwen-14B | $8.77_{\pm3.12}$ / $-16.00_{\pm3.09}$ | $17.41_{\pm3.09}$ / $-9.36_{\pm3.12}$ | $11.58_{\pm3.11}$ / $-12.93_{\pm3.11}$ | $9.56_{\pm3.12}$ / $-17.94_{\pm3.08}$ |
| Qwen-14B-Inst. | $-0.72_{\pm3.13}$ / $-20.50_{\pm3.07}$ | $11.84_{\pm3.11}$ / $-19.22_{\pm3.08}$ | $4.22_{\pm3.13}$ / $-15.59_{\pm3.10}$ | $3.73_{\pm3.13}$ / $-19.79_{\pm3.07}$ |

Table 18: Diff-Bias scores (ambiguous / disambiguated) for the Chinese (CBBQ) and Japanese (JBBQ) datasets, with 95% CI half-widths shown as subscripts. The two datasets are presented as vertical blocks for clarity.

| | w/o RAG | VectorIndex | BM25 | Contriever |
|---|---|---|---|---|
| GPT-3.5 | $5.16_{\pm4.36}$ / $-9.33_{\pm4.35}$ | $11.31_{\pm4.34}$ / $-0.10_{\pm4.37}$ | $17.41_{\pm4.30}$ / $-1.19_{\pm4.36}$ | $9.77_{\pm4.34}$ / $-1.79_{\pm4.36}$ |
| Llama3-8B-Inst. | $5.65_{\pm4.36}$ / $1.59_{\pm4.36}$ | $6.80_{\pm4.36}$ / $-3.97_{\pm4.36}$ | $9.18_{\pm4.35}$ / $-1.88_{\pm4.36}$ | $10.17_{\pm4.34}$ / $-5.06_{\pm4.36}$ |
| Qwen-7B-Inst. | $10.02_{\pm4.34}$ / $-3.67_{\pm4.36}$ | $15.43_{\pm4.31}$ / $-2.08_{\pm4.36}$ | $16.27_{\pm4.31}$ / $-1.39_{\pm4.36}$ | $15.87_{\pm4.31}$ / $-0.10_{\pm4.37}$ |
| Qwen-14B | $3.77_{\pm4.36}$ / $-7.34_{\pm4.35}$ | $0.55_{\pm4.37}$ / $-4.66_{\pm4.36}$ | $7.39_{\pm4.35}$ / $-4.56_{\pm4.36}$ | $5.21_{\pm4.36}$ / $-6.05_{\pm4.36}$ |
| Qwen-14B-Inst. | $-2.38_{\pm4.36}$ / $-2.38_{\pm4.36}$ | $-3.08_{\pm4.36}$ / $-5.95_{\pm4.36}$ | $-0.20_{\pm4.37}$ / $-4.37_{\pm4.36}$ | $-1.64_{\pm4.36}$ / $-4.56_{\pm4.36}$ |

Table 19: Ambiguous / disambiguated Diff-Bias scores (with 95% CI half-widths as subscripts) for different retrieval methods on the full-set.