**S.T. Yau High School Science Award**

**Research Report**

**The Team**

Name of team member: Yijun Wang
School: Tonbridge School
City, Country: Tonbridge, UK

Name of supervising teacher: Ruth Davis
Job Title: Director of Global Futures
School/Institution: Tonbridge School
City, Country: Tonbridge, UK

**Title of Research Report**

Towards Human-Centric Fairness: The VBBQ Benchmark and Sensitive-Attribute Invariance Control for Multimodal Large Language Model Debiasing

**Date**

23.08.2025

# Towards Human-Centric Fairness: The VBBQ Benchmark and Sensitive-Attribute Invariance Control for Multimodal Large Language Model Debiasing

**Yijun Wang**

Tonbridge School

Email: 22wangy@tonbridge-school.org

## Abstract

**Multimodal Large Language Models** (LLMs) already support myriads of applications, and possess formidable multimodal generation and comprehension capabilities. Yet it is precisely due to their extensive pretraining on open web data that they end up perpetuating and even amplifying **social biases and stereotypes**. If we continue to allow sensitive attributes like race or gender to dominate responses and cause discriminatory outputs, there will be grave consequences for LLM applications in critical fields like medicine. Hence, our research is aimed towards building state-of-the-art methods of **mitigating human-centric biases in LLM responses**.

To aid in measuring multimodal biases and alleviate the current lack of **multimodal bias benchmarks**, we extend the existing **Bias Benchmark for QA (BBQ)** dataset for multimodal inputs and compile **vBBQ**, a vision-augmented bias benchmark. vBBQ covers **nine single-dimension biases** (age, disability, gender, nationality, physical appearance, race, religion, socio-economic status and sexual orientation, excluding intersectional categories) and contains a total of **1,092 entity-image pairs**. Image sourcing was based on a hybrid pipeline, containing both synthetic image generation and online searching tools. The final vBBQ set has undergone strict machine and human curation from **more than 100,000 preliminary images**.

In addition to this benchmark, we propose a **Sensitive-Attribute Invariance Control (SAIC)** methodology, and implement it as an agentic framework which allows LLMs to intelligently utilize a range of SAIC tools for sensitive entity removal. Through means such as **textual entity anonymization** and **sensitive visual attribute removal**, our framework makes bias mitigation accessible to end users, unlike traditional bias removal methods such as pre-training data resampling or model training adjustments.

We conducted controlled, large scale experiments to evaluate the effect of our SAIC approach on four LLMs **(GPT-4o, GPT-5, Gemma3 and Llama4)** under four sets of experimental parameters (baseline BBQ, BBQ with SAIC, baseline vBBQ, vBBQ with SAIC). Additionally, adding SAIC directly improves performance in GPT-4o and GPT-5, especially so in previously underexplored social dimensions, and increases LLMs' bias mitigation versatility over many dimensions. Aggregated results show that GPT-5's overall accuracy increases by **+9.5%** for BBQ and **+8.7%** for vBBQ, as well as improvements of **+6.7%** and **+7.6%** for GPT-4o. Models showed stronger gains in disambiguated settings, with a mean **+17.9%** BBQ increase and a mean **+17.1%** vBBQ increase for GPT-5. The religion category proved to be the area of greatest gain, where SAIC improves GPT-5's accuracy performance by **+47.5%** in disambiguated BBQ and **+39.2%** in disambiguated vBBQ.

All our source code (including tools) is made open-source on: https://github.com/jwangexe/vision_bbq. The VisionBBQ dataset is available at Google Drive.

# Acknowledgments

**Commitments on Academic Honesty and Integrity**

We hereby declare that we

1.  are fully committed to the principle of honesty, integrity and fair play throughout the competition.
2.  actually perform the research work ourselves and thus truly understand the content of the work.
3.  observe the common standard of academic integrity adopted by most journals and degree theses.
4.  have declared all the assistance and contribution we have received from any personnel, agency, institution, etc. for the research work.
5.  undertake to avoid getting in touch with assessment panel members in a way that may lead to direct or indirect conflict of interest.
6.  undertake to avoid any interaction with assessment panel members that would undermine the neutrality of the panel member and fairness of the assessment process.
7.  observe the safety regulations of the laboratory(ies) where we conduct the experiment(s), if applicable.
8.  observe all rules and regulations of the competition.
9.  agree that the decision of YHSA is final in all matters related to the competition.

**We understand and agree that failure to honour the above commitments may lead to disqualification from the competition and/or removal of reward, if applicable; that any unethical deeds, if found, will be disclosed to the school principal of team member(s) and relevant parties if deemed necessary; and that the decision of YHSA is final and no appeal will be accepted.**

*(Signatures of full team below)*

*Yijun Wang.*
_____
Name of team member: Yijun Wang

*Ruth Davis*
_____
Name of supervising teacher: Ruth Davis

Declaration of Academic Integrity

The participating team declares that the paper submitted is comprised of original research and results obtained under the guidance of the instructor. To the team's best knowledge, the paper does not contain research results, published or not, from a person who is not a team member, except for the content listed in the references and the acknowledgment. If there is any misinformation, we are willing to take all the related responsibilities.

Names of team members: Yijun Wang

Signatures of team members: *Yijun Wang*

Name of the instructor: Ruth Davis

Signature of the instructor: *Ruth Davis*

Date: 23.08.2025

# Contents

# 1 Introduction

**Large Language Models** (LLM) have come to dominate the field of **Natural Language Processing** (NLP) in recent years. This is due to their impressive ability to extract logical connections from within prompts and produce convincing responses via probabilistic next-token prediction. Even more recently, OpenAI have announced the release of the **GPT-5** model as a "unified system" which "puts expert-level knowledge into everyone's hands". [42] Other examples of notable decoder-only LLMs include **Gemini** [58], **DeepSeek V3** [26] and **Llama 4** [56].

LLMs have been used in a great variety of practical applications to date. For instance, LLMs can be leveraged in the fields of medicine (to make diagnoses [2], formulate treatment plans [9] or communicate with patients [7]) and scientific discovery (to generate novel hypotheses backed by existing literature [64, 70]). The main advantage of LLMs in these fields is their ability to rapidly absorb large quantities textual information in a way which is still challenging to humans.



***Figure 1:*** *Visualized summary of our model evaluation pipeline. The far left shows a multimodal vBBQ example, which is passed to multimodal LLMs first for tool calls, then for a final response.*

However, with great power comes great responsibility. Despite existing advancements in bias mitigation, LLMs still continue to learn and amplify **social biases** present in training data [25]. Hence, it is possible for **discriminatory outcomes** to occur in the aforementioned systems due to bias in the LLMs, which would have serious real-world impacts [30, 46]. For instance, if diagnosis or patient communications systems discriminated by race or gender [23], targeted groups would almost certainly receive worse-quality care. For these reasons, **model bias** is a concerning phenomenon which has several causes, such as **model architecture design choices** and **imbalanced training data**. Yet in order to evaluate the effect of various bias mitigation measures in the first place, **robust bias evaluation metrics** are required, especially in **multimodal settings**.

5

In fact, due to the relative lack of attention towards bias reduction, **evaluation bias** may be introduced during the evaluation phase of LLMs. Many widely used bias benchmarks, such as **BLEU** [44] for machine translation and **MMLU** [16] for general knowledge, are not designed to measure **model bias**, showing that bias evaluation was not previously a priority for LLM trainers. However, in practice, maximum performance on traditional benchmarks without including bias benchmarks **can inadvertently increase bias** in the resulting models, as learning more useful correlations requires also learning more biased correlations. Wu and Aji [68] provide further evidence that current model evaluation systems may be **misguided**, as the paper shows that **both** human and LLM evaluators in the LLM Elo rating system prefer responses with **factual errors** over responses with **grammatical errors**.

Even though many bias benchmarks now exist, such as **BBQ** [45], **WinoBias** [71] and **StereoSet** [35], they focus on **text prompts** rather than **multimodal inputs**, which are under-investigated. One reason for this relative lack of multimodal bias benchmarks is the difficulty of gathering sufficiently relevant images while maintaining real-world authenticity. On one hand, while generated synthetic images are **easier to curate**, they may **neglect the significance of bias** in real-world contexts [37]. On the other hand, images sourced from online can be **lower-quality** or **more challenging to gather**, even if they **better imitate real-world multimodal inputs**. In order to provide a complete framework for evaluating multimodal bias, more benchmarks containing images of both types should be created.

Additionally, many bias mitigation methods are **impractical for end users**. Traditional bias mitigation techniques reduce bias in one of three ways: **preprocessing training data** to remove bias-causing imbalances, **alter model training methods** to minimize the bias absorbed, or **adjust model prompts and outputs** to reduce bias. While **pre-processing** and **intra-processing** methods have achieved success, these require access to **large datasets** and **hardware** capable of model pre-training, which make them generally inaccessible to end users. Therefore, further investigation into currently-unexplored agentic post-processing methods would aid in the propagation of **fair LLM systems** to all end users.

To address these gaps in current multimodal bias research, we introduce the **vBBQ** (short for VisionBBQ) **dataset**, which is designed to **both identify and mitigate biases** in LLMs **when used in conjunction with BBQ**. We also held **large-scale, controlled tests** of the vBBQ benchmark with four LLMs (GPT-4o, GPT-5, Llama4 and Gemma3) as illustrated above in Figure 1, and set out a **Sensitive-Attribute Invariance Control** (SAIC) methodology making use of three **callable tools**, which is designed to **agentically** mitigate multimodal bias. **The vBBQ dataset consists of two parts: images for entities in all base categories, as well as metadata on those entities. The visual component is a hybrid of generated and online-sourced images, while the metadata consists of unique entity identifiers, entity placeholder names, visual trait tags, and the representative image of that entity.**

We generated **exactly one image** and sourced **roughly 100 images** using Google Cloud **per distinct entity**, relying on **visual trait tags** generated by GPT-4.1. Following image acquisition, a **rigorous automatic filtering process** was enacted on online images via GPT-4.1, which defaults to using generated images if none are suitable. Lastly, a random selection of representative images was manually checked for errors or ambiguities.

The **main contributions** of this paper can be summarized as follows.

- Firstly, we created **vBBQ**, a benchmark image dataset designed to augment BBQ text prompts. This dataset contains **1,092 quality-checked images** which augment **all entities** in BBQ except intersectional biases. Each image has undergone a **rigorous curation process**, consisting of **multiple rounds** of human and machine evaluation. We hope that this new dataset, which contains **clear, human-checked images** for all entities in nine bias classes, will expedite **future investiga-**

**tions of multimodal bias** in LLMs.

- As a means to mitigate bias in prompts, this paper also presents a **Sensitive Attribute Invariance Control** (SAIC) approach which allows LLMs to independently utilize **three functional tools** (image grayscaling, image sketching and entity obfuscation). This **agentic** method of bias mitigation, which has been **empirically shown** to reduce the frequency at which LLMs exhibit bias, is also designed to be **accessible** to the average end-user.

- Finally, using vBBQ, we conducted **a large-scale, controlled evaluation** of the extent of bias present in the four LLMs **GPT-4o, GPT-5, Llama4 and Gemma3** in various multimodal settings. These experiments showed that our proposed SAIC approach greatly increases the GPTs' accuracy scores for all categories in disambiguated settings. After integrating a SAIC-based pipeline, the newly-released GPT-5 undergoes the **highest overall accuracy gain**, by **+9.5%** for BBQ and **+8.7%** for vBBQ, with GPT-4o close behind at **+6.7%** for BBQ and **+7.6%** for vBBQ. The *Religion* category proved to be the area of greatest gain, where SAIC improves GPT-5's accuracy performance by **+47.5%** for BBQ and **+39.2%** for vBBQ in disambiguated settings. Gemma3 saw a **+2.6%** gain for BBQ but a **-0.7%** loss for vBBQ. The two GPT models exhibited much stronger gains in disambiguated settings, with a mean **+17.9%** BBQ increase and a mean **+17.1%** vBBQ increase for GPT-5.

The rest of the paper is organized as follows. **Section 2** reviews the extent to which bias is present in LLMs, and explores existing bias benchmarking and mitigation methods in detail. In **Section 3** we detail our overall methodology, with **Section 3.1** covering the curation of the vBBQ dataset and **Section 3.2** pertaining to the tools used during the model prediction stage. **Section 4** describes the precise experimental environment and setup, details the various bias metrics used to make comparisons, and gives the results achieved from the LLM testing. **Section 5** concludes the paper by interpreting these results and suggesting next steps in the field of multimodal bias benchmarking.

# 2 Related Work

## 2.1 The Evolution of Large Language Models

**Foundation Era (2017-2020)**



*Figure 2: Timeline of major language model developments color-coded by architecture type*

**Large Language Models** (LLMs) have become a central topic of research in the field of **natural language processing** (NLP), due to innovations in **neural network** (NN) architecture, pretraining on large text corpora and improved computer hardware capabilities. Prior to the **Transformer architecture**, early attempts at NLP using deep learning relied on recurrence through **Recurrent Neural Network** (RNN) or **Long-Short Term Memory** (LSTM) approaches. [17, 33, 57] These models processed text sequentially, which limited parallelization and restricted the scale of training datasets.

**Transformers** replace recurrence with highly parallelized **self-attention mechanisms**, which increases model efficiency when used in conjunction with modern GPUs. [61] This key innovation in neural network architecture paved the way for two main pre-training approaches. **Masked language modeling** (MLM) masks random tokens in text during training, which allows LLMs to make use of context from tokens before and after those unknown tokens. [11] During training with **causal language modeling** (CLM), LLMs must predict unknown tokens from left to right, so can only use context from already generated tokens. [48] While MLM is the basis for encoder-only models such as **BERT** [11] and **RoBERTa** [28], CLM is used by decoder-only models such as **GPTs**. **Encoder-only models** usually possess strong text comprehension capabilities and are well-suited to question answering tasks. [11] On the other hand, **decoder-only models** tend to be better at text generation than encoder-only models [48, 49]. Table 1 below lays out the differences in architecture, training method and characteristics of encoder-only and decoder-only LLMs.

*Table 1: Table showing the differing architecture and characteristics of encoder-only and decoder-only LLMs.*

| LLM Architecture | Pre-Training Technique | Strengths | Examples |
|---|---|---|---|
| Encoder-only | Masked language modeling | Comprehension, QA | BERT, RoBERTa |
| Decoder-only | Causal language modeling | Text generation | GPTs |

As shown in Figure 2, OpenAI released new versions of GPT-3 (now known as GPT-3.5) in 2022 with the ability to insert text or edit existing text [4]. This feature led to the use of Codex, a fine-

tuned model built on GPT-3.5, for the code completion tool Github Copilot [24]. ChatGPT was also a fine-tuned model built on GPT-3.5, which underwent Supervised Fine-Tuning and Human Feedback Reinforcement Learning during the fine-tuning process. [39]

Released in 2023, GPT-4's NLP capabilities surpassed those of previous models when evaluated on a diverse set of benchmarks, including exams originally designed for humans. Notably, the model achieved a score in the 90th percentile in the Uniform Bar Exam, whereas GPT-3.5 only achieved a score in the 10th percentile. GPT-4 was also the first multimodal GPT model, allowing it to accept images as input and complete specified vision tasks [40]. GPT-4o, released in 2024, surpassed GPT-4's original multimodal capabilities after undergoing training from multimodal data. The model accepts any combination of text, image, audio or video and passes outputs as text, image, and/or audio. [41]

## 2.2 The Manifestations and Consequences of Bias in Large Language Models

LLMs have consistently demonstrated strong capabilities in text generation, comprehension, and multimodal tasks. However, these models will inevitably contain biases, which can have serious implications for the real-world applications of LLMs. In this paper, we define bias as an unfair tendency in a model which leads to discriminatory output. [51]

As described previously, decoder-only LLMs such as GPT are trained to predict the next token in a large corpus of human-generated text. As a result, these models are able to perpetuate existing social biases within training data [5, 8]. As demonstrated by [5], social biases can be found within word embeddings, which are an essential first step in the LLM pipeline from which biases may propagate through the whole model. Word embeddings map words to numerical vector representations, with the special property that vector similarity correlates with semantic similarity [31]. For the w2vNEWS word embedding trained on Google News articles, [5] showed that the difference between the embeddings for "man" and "woman" is nearly that of "programmer" and "homemaker", like so: [5]

$$\overrightarrow{man} - \overrightarrow{woman} \approx \overrightarrow{programmer} - \overrightarrow{homemaker} \tag{1}$$

, which implies that there is a distinction between typical male and female jobs within the embeddings. These biases can lead to gender stereotypes forming, which adversely affect model responses.

Kotek et al. [23] corroborate these results over the entirety of the LLM, by testing on a dataset inspired by the WinoBias set. Other studies have also shown that LLMs exhibit bias in many other social dimensions, including (but not restricted to) religion, occupation, political views, sexual orientation, nationality, and disability. [1, 43, 63, 62, 12]

## 2.3 Established Bias Mitigation Strategies

Due to bias' harmful implications for large language models and the ethical concerns which arise as a result, measures which reduce bias in model outputs are essential for the continued use of LLMs.

Bias mitigation techniques generally fall into one of three categories: pre-processing, intra-processing and post-processing, depending on the stage at which intervention occurs. In most cases pre-processing techniques debias model training data, intra-processing techniques augment the model training process, and post-processing techniques adjust LLM outputs.

### 2.3.1 Debiasing Training Data via Balancing Methods

Pre-processing techniques typically aim to make the LLM's training data more balanced in an attempt to prevent biases from emerging during training. Examples of data balancing strategies include

oversampling from underrepresented groups, undersampling from overrepresented groups, and synthetic text generation. [14] One widely-used oversampling technique is Counterfactual Data Augmentation (CDA) proposed by [29]. In CDA, generic gendered words in training text such as "woman" are inverted (i.e. to "man"), and the resulting text is added back to the dataset. For example, suppose that this sentence is found in the corpus:

*"The police officer succeeded because he was strong and intelligent."*

CDA would replace "he" with "she" and add the following sentence to the dataset:

*"The police officer succeeded because she was strong and intelligent."*

Note that both sentences would be present in the final balanced dataset. However, since CDA avoids inverting first names and pronouns associated with them, sentences such as these would be duplicated and might amplify gender bias:

*"Bruce succeeded because he was strong and intelligent."*

[32]'s proposal of Counterfactual Data Substitution (CDS) resolves this issue with first names using bipartite graph matching. In CDS, *Bruce* might be replaced by *Irene*, as they tend to represent opposite genders and have similar frequency. Hence CDS would return the sentence:

*"Irene succeeded because she was strong and intelligent."*

### 2.3.2 Innovations in the Model Training Process

Intra-processing methods tend to be applicable to already pre-trained models as part of a fine-tuning process, or involve training additional models within a debiasing pipeline. For instance, soft-prompt tuning is a fine-tuning technique where a series of token embedding vectors $[V_1, V_2, V_3, \dots]$ is appended to input data. These vector "soft prompts" are then trained, while the original model parameters are frozen. Unlike hard-prompt tuning which uses specific tokens as prompts, these embedding vectors can encode features which do not exactly match with any token.

[59] used this technique to preserve model bias from pre-training data for more accurate evaluation and comparison, while [13] utilizes it to fine-tune BERT and RoBERTa models for mitigation of gender bias without the risk of catastrophic forgetting. [34] proposes a Semi-Parametric Editing with a Retrieval-Augmented Counterfactual Model (SERAC) method, which allows model editors to inject updated knowledge or correct behaviors into LLMs.

Firstly, "edits" to the model's knowledge are stored in a separate "edit memory". Each time the model is prompted, a scope classifier determines whether the question requires information from that memory. If so, a new counterfactual model overrides the original model to process the prompt. This approach does not require base model retraining each time new edits are added, making it suitable for small adjustments.

As described in [15], LLMs such as BERT can be further pre-trained on a relatively small corpus with the aim of bias mitigation. To achieve this result, the researchers introduced an equalizing loss function to equalize the association between gender-neutral words such as "intelligent" and gender-specific terms like "man" or "woman". Furthermore, they disrupted clusters of words forming, such as *entrepreneurs, protégé, aspiring, arrogant, bodyguard*, by pre-training with a special declustering loss function.

### 2.3.3 Post-processing Calibration for Unbiased Outputs

Post-processing methods occur after model training is complete. Some approaches utilize prompt engineering, while others alter model output in a way which reduces bias. One of example of this is Kaneko et al [21], which showed that LLMs achieved lower bias scores when evaluated on gender bias with chain-of-thought (CoT) prompting. CoT is a prompt engineering technique which encourages LLMs to generate a series of intermediate reasoning steps before giving their final answer. [66]

The Decoding-time Experts (DExperts) method, proposed by in 2021 by Liu et al [27], relies on logits predicted by non-biased "expert" models $LLM_{+e}$ and biased "anti-expert" models $LLM_{-e}$ for the original prompt $p$.

Initially, all three models (including original model $LLM$) give their predictions as logits $z, z^+, z^-$. Then a final logit $z^{debiased}$ is computed with a constant $\alpha$ like so:

$$z^{debiased} = z + \left(\sum z^+ - \sum z^-\right) \times \alpha \tag{2}$$

, such that tokens are only chosen if they are considered likely by the experts and unlikely by the anti-experts.

## 2.4 Existing Bias Evaluation Frameworks

Researchers have created several metrics and datasets to measure the bias present in model outputs. These tools aid in the identification of areas where the model is particularly biased, providing direction to mitigation efforts. They can also act as benchmarks which measure the effectiveness of bias mitigation.

### 2.4.1 Bias Benchmark for QA

The original Bias Benchmark for QA (BBQ) dataset was proposed in 2022 by Parrish et al. [45] to benchmark LLM biases for the task of question answering. The paper which presented this dataset also used the data to determine how biased model responses are (in an ambiguous question context), and if social biases drive LLMs to override factually correct answers (in a disambiguated context containing a clear, correct response).

The dataset covers 9 social dimensions: age, disability status, gender identity, nationality, physical appearance, race/ethnicity, religion, socio-economic status (SES) and sexual orientation. The researchers used a total of 325 distinct handwritten templates which, after substituting entity names, can generate 58,492 unique examples. Examples from each template are validated via Amazon Mechanical Turk, and only templates which achieve a minimum of 4 gold labels out of 5 annotators are included in the dataset.

Each template in BBQ contains ambiguous context and disambiguating context. The difference between the two is that while disambiguating context gives strong evidence for one entity being the answer, the ambiguous context introduces the scenario and does not give sufficient information on either entity for a definite answer. Templates also contain one positive and negative question each, to determine whether models truly associate disadvantaged groups with negative stereotypes, or if they are simply more likely to give certain groups as responses. Note that the term "negative question" in this context refers to a violation of general US social values, with the "positive question" being its converse. (i.e. "Who was forgetful?" is the negative question and "Who was not forgetful?" is the positive question)

In order to determine the model bias in neutral settings, the researchers used ambiguous questions without disambiguating context to avoid favoring one entity over the other. In this case, the only unbiased answer which the model can give is the one labeled `unknown` (or equivalent). Then, to determine if the model missed clearly correct responses due to bias, the researchers used disambiguated questions which

contain both ambiguous and disambiguating context. Here, the unbiased answer is determined by the exact template used and the question polarity.

### 2.4.2 Extensions of the BBQ Framework

The original Bias Benchmark for QA dataset has inspired a number of other benchmark datasets which aim to evaluate and quantify bias present in LLMs. These variants build upon BBQ's framework for bias evaluation and apply it to other languages and task formats. Among the languages which BBQ has been adapted to are Chinese [18], Korean [19], Dutch, Spanish, Turkish ([38] for all three), Japanese [69], German [54], and Basque [53].

Jin et al [20] also proposed the Bias Benchmark for Generation (BBG) dataset, which applies BBQ's concept to measure bias in LLM text generation. BBQ's disambiguating context is first anonymized by replacing entity names with ambiguous labels "one" and "the other". Then the evaluated LLM is prompted to generate a continuation of the scenario using both ambiguous and disambiguating context, to ensure the relevance of the continuation. The researchers combined the seed story, model continuation and BBQ questions for use as a machine reading comprehension (MRC) passage. Their use of GPT-4 as an evaluator model allowed them to determine whether a generation kept the characters ambiguous (unbiased), aligned with stereotypes (biased), or opposed stereotypes (counter-biased). Table 2 below contains detailed information on the aforementioned BBQ variant datasets, such as their expected response format, year of publication, language, template source, and so on.

***Table 2:*** *Comparison of BBQ dataset variants and related benchmarks. The table shows key characteristics including response format, language coverage, template sources, scale, and annotation methodology.*

| Name | Response Format | Year | Language | Template Source | Num. Templates | Num. Social. Dims | Num. Examples | Annotation Method |
|------|-----------------|------|----------|-----------------|----------------|-------------------|---------------|-------------------|
| BBQ | Multiple-choice (MC) | 2022 | English | Experts | 325 | 9 | 58,492 | Crowdsourced |
| CBBQ | MC | 2023 | Chinese | Online | 3,039 | 14 | 106,588 | Experts |
| KoBBQ | MC | 2024 | Korean | BBQ | 268 | 12 | 76,048 | Crowdsourced |
| MBBQ | MC | 2024 | Dutch, Spanish, Turkish | BBQ | 98 | 6 | 10,072 | Experts |
| JBBQ | MC | 2025 | Japanese | BBQ | 245 | 5 | 50,856 | Experts |
| GG-BBQ | MC | 2025 | German | BBQ | 50 | 1 (Gender) | 5,936 | Experts |
| BasqBBQ | MC | 2025 | Basque | BBQ | 215 | 8 | 43,240 | Experts |
| BBG | Open-Ended + MC | 2025 | English, Korean | BBQ/KoBBQ | 232/286 | 9/12 | 82,136/38,316 | Experts |

### 2.4.3 Other Bias Evaluation Benchmarks

CrowS-Pairs (Crowdsourced Stereotype Pairs, 2020) [36] is an English-language benchmark containing 1508 sentence pairs spanning 9 social dimensions (race, gender, sexual orientation, religion, age, nationality, disability, physical appearance, SES). In CrowS-Pairs, each sentence pair contains a more biased and a less biased sentence, which models are asked to choose from. The dataset was created via crowdsourcing on Amazon Mechanical Turk.

StereoSet (2020) [35] is an English-language benchmark containing 16,995 Context Associated Tests (CAT) spanning 4 social dimensions (gender, profession, race and religion). CATs contain an unclear context (similar to ambiguous context in BBQ) and 3 options for an LLM to choose from. One choice always confirms a stereotype, one always opposes a stereotype, and the last one is always unrelated to the question. For instance, the context

<div align="center">

*Girls tend to be more _____ than boys*

</div>

matches the options *soft* (stereotype), *determined* (anti-stereotype) and *fish* (unrelated). The dataset was created via crowdsourcing.

EEC (Equity Evaluation Corpus, 2020) [22] is an English-language benchmark containing 8640 sentences spanning 2 social dimensions (race and gender). During evaluation, LLMs give each sentence-emotion pair $(S, A)$ a real-valued sentiment intensity score between 0 and 1. The dataset was created via 11 templates, but the original paper does not include a human performance baseline. WinoBias (2020) [71] is an English-language benchmark covering gender bias, which contains 3160 sentences. The dataset was created via two basic sentence structures with crowd-sourcing.

GAP (Gendered Ambiguous Pronouns, 2020) [65] is an English-language benchmark containing 8908 ambiguous pronoun-name pairs concerning gender. These example text sequences contain two nouns and an ambiguous pronoun, which models have to assign to a noun. Pronouns can be placed before both nouns (InitialPro), between them (MedialPro) or after them (Final Pro). For instance, the text

Judging that **he** is suitable to join the team, *Butcher* injects *Hughie* with a specially formulated mix.

is an example from the InitialPro group. The dataset was created via crowdsourced annotation of examples from Wikipedia, and the paper includes several baselines based on heuristic approaches.

BiLG (Bias in Language Generation, 2020) [55] is an English-language benchmark containing 302 prefix templates spanning 3 social dimensions (race, gender and sexual orientation). Each prefix template is a brief prompt such as *The woman worked as...*, which OpenAI's medium-sized GPT-2 model completed (i.e. *...a prostitute under the name of Hariya*). The dataset was created via a combination of researcher-chosen templates, GPT-2 generation and crowdsourced annotation of examples. IndiBias (2020) [52] is a combined English and Hindi-language benchmark containing 1600 sentence pairs which span 7 social dimensions (gender, religion, caste, age, region, physical appearance and SES). Three intersectional categories (gender-religion, gender-caste, gender-age) are also included. The dataset was created via translation from CrowS-Pairs. Table 3 below compares these benchmarks, showing their basic information such as year published, language, expected response type, number of social dimensions, example count and annotation method.

**Table 3:** *A comparison of other bias benchmarks, containing information on their year of publication, language, response type, number of social dimensions, number of examples and annotation method.*

| Name | Year | Language | Response Type | Social Dim. | Examples | Annotation |
|------|------|----------|---------------|-------------|----------|------------|
| CrowS-Pairs | 2020 | English | MC | 9 | 1,508 | Crowdsourced |
| StereoSet | 2020 | English | MC | 4 | 16,995 | Crowdsourced |
| EEC | 2020 | English | Open-Ended | 2 | 8,640 | Crowdsourced |
| GAP | 2020 | English | MC | 1 | 8,908 | Crowdsourced |
| BiLG | 2020 | English | Open-Ended | 3 | 302 | Crowdsourced |
| IndiBias | 2020 | English, Hindi | MC | 7 | 1,600 | Experts |

# 3 Methodology

## 3.1 Constructing a Multimodal Bias Benchmark



**Figure 3:** *Visual overview of the vBBQ dataset creation process, from reading textual BBQ questions to generating entity-image pairs.*

The curation of the VisionBBQ (vBBQ) multimodal bias benchmarking dataset, which addresses the current lack of comprehensive multimodal bias benchmarks, is one of the primary contributions of this paper. As shown in Figure 3, we used a hybrid web searching-generation pipeline to compile 1,092 entity-image pairs from the original BBQ dataset, which was sufficient to provide multimodal context for 31,372 BBQ-style multiple-choice questions. vBBQ dataset curation also took place alongside a strict quality control process where multiple rounds of machine and human review took place, in order to ensure the high quality of all multimodal context gathered.

### 3.1.1 VisionBBQ: A Comprehensive Structural Framework

The vBBQ dataset consists of image and metadata files organized as separate directories within a root directory, and is intended for use in conjunction with the original BBQ dataset. For each entity present in the BBQ data (except for those pertaining to intersectional bias), vBBQ contains roughly 100 images with online sources, as well as 1 AI generated image.

To permit the efficient re-use of images for questions containing previously seen entities, we structured vBBQ as a collection of entity-image pairs. Metadata on all entities is stored in the *dictionary/* subdirectory. For each bias category, *dictionary/* contains a CSV file with relative path *dictionary/{bias_class}_entity.csv*. The searched images of each entity are stored in the *images/{entity_name}/* directory (with each distinct entity occupying its own sub-directory), while the AI generated image for each entity is stored at path *ai_images/{entity_name}.jpg* (where each entity is represented by a single image file). Figure 4 visualizes this file hierarchy more clearly, and gives examples from each dataset location.

***Figure 4:*** *vBBQ file structure diagram showing CSV metadata, generated and searched image directories, including one example from each.*

### 3.1.2 Extracting Entity Visual Attributes

Before compiling vBBQ images, it was necessary to extract a set of entities $E$ from the BBQ question set $Q$, with each entity $ent \in E$ containing metadata required for image sourcing (e.g. visible feature tags). Each of the CSV metadata files present in *dictionary/* contains these columns for all entities $ent$[1]. Firstly, `bbq_id` (or $id_{ent}$) is a unique numerical identifier for each entity which distinguishes entities by bias class, row number in the original BBQ set and label number. For each entity $ent$ on row $i_{row}$ of the bias category uniquely marked $i_{file}$, its identifier $id_{ent}$ is thus calculated:

$$id_{ent} = k_1 k_2 i_{file} + k_2 i_{row} + i_{ent} \tag{3}$$

where $k_1 = 1000000, k_2 = 10$ are constants and $0 \leq i_{ent} \leq 2$ is the entity's answer label (i.e. whether it occupied the A, B or C choice). Although our methodology does not make use of unique entity IDs, we expect that the ability to differentiate between two entities with the same name will be useful for future experiments, especially if researchers plan to obtain images using the full prompt as context.

`name` (or $name_{ent}$) is an verbal identifier sufficient to characterize the entity in the context of a BBQ question (e.g. `The grandfather`). Therefore, we have defined most entity names as the answer text for that entity's corresponding test case. However, in the *Race/Ethnicity* bias category the BBQ researchers [45] gave human names to many entities (e.g. *Fatima al-Fasi*), generated from a list of popular first and last names for each race-gender combination (with data from US censuses). For each answer choice $c$ in the BBQ question set, we split it into its component words $(w_1, w_2, ..., w_n)$ by spaces, and then determined whether any components were present in the first and last name lists given by BBQ, here

---

[1]Note that `mono-spaced text` denotes column names in the CSV files, while *mathematical* formatting denotes the notation used in the rest of this paper.

denoted by $N$. In other words,

$$isName(c) = \begin{cases} 1 & \text{if } \exists w \in c \text{ such that } w \in N \\ 0 & \text{otherwise} \end{cases} \tag{4}$$

where $isName(c)$ denotes whether $c$ is a name introduced by BBQ researchers.

All names within the set of entity names were then replaced by their corresponding race and gender (so *Fatima al-Fasi* would become *F-Arab*), using metadata given within BBQ. This is because the common names used in BBQ rarely convey significant information outside of gender and race, which also happen to be two visible characteristics well suited to image representation.

`tags` (or $T_{ent}$) is a string of comma-separated visible traits which give a reasonably accurate visual description of some entity *ent*. In other words, $T_{ent}$ can be understood as an ordered sequence $(t_1, t_2, t_3, \ldots, t_n)$ where each individual tag $t$ is a phrase of length 1-3 words. For instance, the entity name *The grandfather* has tags *"elderly man,gray hair,wrinkles,glasses,kind smile"*.

We used OpenAI's `gpt-4.1-2025-04-14` model, which was prompted to generate 3 to 5 of the most important and visible traits of each entity as tags. These tags were then used as prompts in both image generation and image search. Furthermore, the resulting images' clarity and accuracy to the entity was evaluated by another LLM using said entity tags. Tag generation is used in our methodology because image evaluation using the original entity name is prone to rejection by the LLM. (insert figure) While the model views the original entity `The black man` as a racial label, it views similar tags such as `black skin` as a visual cue, which greatly decreases the likelihood of the model refusing.

`imgpath` represents the path of an image which matches the entity tags well. Specifically, it is the path of the first image which receives a `10/10` in the Image Evaluation phase, or the generated image for that entity (if there are no suitable searched images). See "Image Evaluation" for more details on this process.

### 3.1.3 Multifaceted Visual Data Sourcing

For each target entity *ent* in the metadata entities set, we automatically retrieved around 100 stock photos $S_{ent} = \{I_{ent,1}, I_{ent,2}, I_{ent,3}, \ldots\}$ using the Google Custom Search JSON API *Search*, and generated one $1024 \times 1024$ image $g_{ent}$ using Google's Imagen-4.0 model *Imagen*. Both prompts used were concatenations of the entity's name and tags, with the search prompt containing additional base tags $T_{base}$ such as *"human"* to ensure the relevance of resulting images, as shown below:

$$S_{ent} = Search(concat(name_{ent}, T_{ent}, T_{base}), 100) \tag{5}$$

$$g_{ent} = Imagen(concat(name_{ent}, T_{ent}), 1) \tag{6}$$

After obtaining the searched image URLs from the Google API in batches of 10, each URL was validated via HTTP access within a set timeout to ensure that all images were accessible. The content of these URLs was then saved in the `images/<entity_name>/` directory with numerical filenames for each image. The result of the Imagen model was saved at `ai_images/<entity_name>.jpg`. At least 100 searched URLs were obtained for each target entity, while one image was generated per entity. This is due to the Imagen model's superior ability to tailor the result specifically to the prompt, as opposed to the risk of mismatch between the web-searched images and the entity tags. Thus, many more searched images are required to increase the likelihood of obtaining images suitable for the entity.

On the other hand, since the web-searched images are mostly stock photos sourced from real-world settings, they are likely to be more authentic and represent entities in a more unbiased way than generated images. Therefore, the latter was utilized as a fallback option and were only used if no image matched the prompt to a sufficient degree.

### 3.1.4 Image Curation and Quality Control Protocols

In order to mitigate the potential mismatch between searched images and the entity tags, we evaluated the quality of our images both manually and using LLMs. For each target entity *ent*, we retrieved the set of searched images $S_{ent}$ and its generated image $g_{ent}$, and for each image $s \in S_{ent}$ assigned an integer quality score $0 \leq score(s) \leq 10$. The images are evaluated by OpenAI's mutlimodal GPT-4.1 model, based on clarity and fulfillment of the tags $T_{ent}$, from a scale of 0 to 10 (with 10 being the highest). The system prompt also penalizes discrepancies in number (i.e. a single person will be rejected if the tags contain "group"). Overall, the scoring process for some image $s \in S_{ent}$ can be considered just as in the following equation:

$$score(s) = LLM(s, T_{ent}) \tag{7}$$

where *LLM* denotes the model, which is used as a function.

The evaluation of each entity ent results in one selected image $I_{ent}$. If there exists an image *s* with $score(s) = 10$, that image is chosen to represent the entity in question. Otherwise, since none of the searched images are sufficient, the generated image $g_{ent}$ is used, like so:

$$I_{ent} = \begin{cases} s & \text{if } \exists s \in S_{ent} \text{ such that } score(s) = 10, \\ g_{ent} & \text{otherwise} \end{cases} \tag{8}$$

Finally, all chosen images are manually checked to ensure that they fulfill the criteria given above. The following pseudocode describes the process of image evaluation in greater detail.

---

**Algorithm 1** vBBQ Image Quality Evaluation Protocol

---

1: **Import** Pandas, OpenAI dependencies
2: **Import** modules and entity set $E$
3: $I \leftarrow \emptyset$
4: **for** each entity $ent \in E$ **do**
5:      Load searched image set $S_{ent}$ and generated image $g_{ent}$
6:      $done \leftarrow false$
7:      **for** each image $s \in S_{ent}$ **do**
8:          $score(s) \leftarrow LLM(s, T_{ent})$
9:          **if** $score(s) = 10$ **then**
10:              $I_{ent} \leftarrow s$
11:              $done \leftarrow true$
12:              **break**
13:          **end if**
14:      **end for**
15:      **if** $done = false$ **then**
16:          $I_{ent} \leftarrow g_{ent}$
17:      **end if**
18: **end for**
19: **Return** set of selected images $I$

---

## 3.2 SAIC: An Agentic Framework for Multimodal LLM Bias Mitigation

### 3.2.1 Overview of Our Approach

This paper proposes an agentic method of bias mitigation, which allows the model to independently decide the extent of transformations on the input context. Hence, during the prediction process the model may be given access to a set of functions $U$. Set $U$ may contain the grayscale conversion tool $u_g$, the sketch extraction tool $u_s$ and the entity anonymization tool $u_a$, though some tools may be omitted depending on specific experimental parameters. [2]

If the two images are to be used, they are encoded using base64 before the prediction process for the sake of compatibility with the Ollama API's input format. Additionally, each image was resized such that its longest side measured no more than 512 pixels before the model predictions [3]. An initial request $r_0$ is sent to the model, which yields several structured tool calls $(u, p)$ returned by the API (which call on tool $u$ with parameters $p$), like:

$$LLM(r_0) = \{(u_0, p_0), (u_1, p_1), (u_2, p_2), \ldots, (u_n, p_n)\} \tag{9}$$

Then the requested tools are executed locally and their results override the existing multimodal context to create a new context $r_1$. The model then gives an answer of three choices (based on the BBQ format) based on the updated context $r_1$, which can be expressed as such:

$$LLM(r_1) \in \{\texttt{"A"}, \texttt{"B"}, \texttt{"C"}\} \tag{10}$$

---

[2]Note that the first two are image tools while the third is a text tool.
[3]Note that the images publicly available in the BBQ dataset are full-size. This measure prevents larger images from potentially introducing more cues for biased judgments by maintaining consistent image dimensions

All requests are handled using structured generation to ensure that the output format is always as expected.

### 3.2.2 Removing Chromatic Information via Grayscale Conversion



**Figure 5:** *Illustrated grayscale conversion of an example image: the luminance formula is applied to each pixel of a color image and maps to a new grayscale image.*

The grayscale conversion tool $u_g$ takes a color image $I_c$ as input and converts it into grayscale image $I_g$, following the ITU-R BT.601 standard. Specifically, a luminance value $L$ is calculated for each pixel within the image using the Python PIL library, based on the $R, G, B$ intensities of said pixel (where 0 is black), like [4]:

$$L = 0.299R + 0.587G + 0.114B \tag{11}$$

The process described above is illustrated by Figure 5 using an example image.

Within images, color can often be a significant medium for bias. For example, skin color is a characteristic of ethnicity, while color grading can hint at the image's intended location. As a result, grayscale conversion is a useful transformation to reduce bias since it greatly decreases color information in the image and removes color cues which could misguide models.

### 3.2.3 Retaining Structural Elements via Sketch Conversion

The sketch conversion tool $u_s$ can be treated as an extension of grayscale conversion, which outputs a transformed image $I_s$ containing only the outlines of objects. It is more extreme than grayscale conversion, which retains the shading of depicted objects. In the process of sketch conversion, the image must first undergo the grayscale conversion process described above. Next, each pixel in the image is inverted, so that light areas become dark and vice versa. In practice, this is implemented like so:

$$I = 256 - L \tag{12}$$

where $I$ refers to the inverted pixel and $L$ refers to the original grayscale.

A Gaussian blur with radius 10 is then applied to image $I$, which removes sharp edges and detailed features. This step prepares the inverted image for its use as an overlay by removing sharp features which might otherwise be present in the final sketch. To obtain the final sketch, we applied a dodge blend with

---

[4]Given that unsigned 8-bit integers are used, the range of $L$ is the same as that of $R, G, B$.

*Figure 6: An illustrated example of the sketch conversion pipeline. After preliminary grayscaling, the image is inverted and a Gaussian blur is applied to create a filter. Then a dodge blend between the original image and filter results in a sketch-like image.*

the original grayscale image as the base image, and the blurred, inverted image as the foreground, with the following formula:

$$p_{out} = \min(\frac{b \times 255}{256 - f + 1}, 255) \tag{13}$$

where $b$ represents the base image (here the original grayscale image $L$), and $f$ represents the foreground image (here the transformed image $I$). Figure 6 depicts part of this process, starting from a grayscaled example image and ending up with a sketch.

The process of sketch transformation only preserves the basic outlines of objects in the image, which removes potential biases conveyed by the image's grayscale shade or lighting. In this respect it is a more extreme image-abstraction method than simply grayscaling.

### 3.2.4 Eliminating Bias Targets using Entity Anonymization

Entity anonymization $u_a$ is a textual debiasing function tool where the entities $(ent_1, ent_2)$ present in a prompt $p$ are replaced with generic names such as `Entity1` and `Entity2`, as shown in Figure 7. This removes information about the entities present in textual prompts, which in turn prevents the LLM from making biased judgments based on those entities' characteristics.

In order to ensure the validity of the original BBQ labels, an LLM assistant was prompted to output a supplementary "answer key" of answer choices to entities $K$, and to maintain consistency in the replaced entity names. The overall tool can be represented like so:[5]

$$u_a(p_0) = (p_1, K) \tag{14}$$

$$K = \{A \mapsto x,\ B \mapsto y,\ C \mapsto z\} \tag{15}$$

We implemented this tool using OpenAI's GPT-4.1 model, however any LLM with structured output capabilities would be functionally equivalent.

---

[5]Two of $(x, y, z)$ are in $(ent_1, ent_2)$ and the last one is the unknown response $ent_?$

*Figure 7: Informative diagram showing the function of the entity anonymization tool.*

Unlike rule-based parsing approaches, LLMs have the ability to parse text flexibly and adapt to the unpredictable format of real-world prompts. For instance, rewording entity names ("The old man" to "The elderly man") could cause rule-based approaches to fail or become impractical, but would not cause much additional difficulty for LLMs.

### 3.2.5 Improving Response Consistency with Structured Output

During predictions, we forced the model response to adhere strictly to a defined JSON schema to ensure response validity. For instance, the predicted label was required to be a string out of `"A"`,`"B"`,`"C"`. The precise implementation of Structured Output requires a technique known as Constrained Decoding [47]. At each step in this process, model generation is restricted to tokens which could be valid if added to the current output. For instance, if the model was required to generate the schema `{"name":` `"Harry"}`, after it has already generated `{"nam` it would be restricted from choosing invalid tokens such as `"]` which would yield invalid JSON.

This technique is sometimes implemented by converting a given schema into a context-free grammar(CFG) [47], which is a set of rules defining the structure of the response. At each step of generation the model computes a set of valid tokens from existing content and CFG rules efficiently, using a cached data structure. However, other implementations, such as with Finite State Machines, also exist [67]. In practice, we utilized the OpenAI API's structured output feature (which uses CFGs) and the Pydantic library for the strict typing of response fields.

### 3.2.6 Tool Calling: Foundation of the Agentic Approach

Although LLMs have a wide range of textual capabilities, they are generally unable to transform images, consult APIs or execute code by themselves. Therefore, in order to mitigate biased inputs via an agentic method, LLMs must be able to reliably call upon external tools (such as those in our toolbox).

To make use of tool calling, we needed to provide a schema $F_u$ for each tool $u$ to the LLM, which contains their name, parameter-and-type pairs $(p, \tau)$, and a brief description of the tool provided. The schema can be represented as the following mathematical object:

$$F_u = (name_u, \{(p_1, \tau_1), (p_2, \tau_2), (p_3, \tau_3), \ldots, (p_n, \tau_n)\}, desc_u) \tag{16}$$

This step informs the LLM of which tools are available for their use, the scope of their functions, as well

as the precise parameter requirements of each.

After the input prompt $p$ is given, the model evaluates which tools are needed for the task and returns tool calls in a specific JSON format (containing the tool name and its argument values) [24], in the following abstracted way:

$$LLM(p) = \{(name_{u_1}, \pi_{u_1}), (name_{u_2}, \pi_{u_2}), (name_{u_3}, \pi_{u_3}), \ldots, (name_{u_n}, \pi_{u_n}),\} \tag{17}$$

where $\pi$ denotes the collection of parameters used. Note that returning a structured JSON object requires the Structured Output feature mentioned earlier. The client then processes these tool requests and runs each tool locally, and includes tool results as context for another request. Now the LLM uses those results to generate the final response. As demonstrated by the experiments we conducted, a range of proprietary and open-source LLMs demonstrate support for structured output and tool calling, so the methodology presented in this paper is model-agnostic.

# 4 Experiments

## 4.1 Experimental Environment

Dataset curation and initial GPT-4o evaluation computations were performed via API calls to the OpenAI (tag generation, quality filtering, testing) and Google Cloud (cloud storage, image sourcing, image generation) platforms. Consequently, no high-end hardware was required for these except for a standard machine (Apple M2, integrated GPU, 16GB RAM). This detail highlights the superior practicality of our post-processing debiasing method made by users, for users.

For testing on open-source models such as Gemma3, we used the Ollama platform running on a Ubuntu 24.04 Linux server with 2 Nvidia GeForce RTX 4090 GPUs and 252GB of RAM. All code was run on a Python 3.12.7 virtual environment, with the following notable third-party libraries used: the OpenAI and Google Cloud client-side APIs, Pandas (data processing) and Pillow (image processing).

## 4.2 Model Selection

We evaluated the performance of several LLMs on the BBQ and vBBQ benchmarks to measure their ability to draw out multimodal bias, as well as the effectiveness of our proposed toolkit to mitigate such bias. Alongside OpenAI's GPT-4o and GPT-5 models, two smaller, open-source models (Gemma 3 and Llama 4) were evaluated using the Ollama platform. All models tested were shown to have vision and tool-calling capabilities.

*Table 4: A list of models evaluated with the BBQ and vBBQ benchmarks.[6]*

| Model Name | Parameter Count | Creator | Open Source? |
|---|---|---|---|
| GPT-4o | ≈200B | OpenAI | No |
| GPT-5 | Unknown | OpenAI | No |
| Llama4 | 109B | Meta | Yes |
| Gemma3 | 12.2B | Google | Yes |

As shown in Table 4, the model parameter counts range from a minimum of 12.2 billion to a known maximum of $\approx 200$ billion, so that we can compare the benchmark's effectiveness on models of varying sizes.

## 4.3 Investigating Bias Across Modalities

Due to vBBQ's nature as an extension to the BBQ benchmark, the LLM evaluation methodology presented in this paper is similar to that proposed in Parrish et al [45]. Both BBQ and vBBQ are QA benchmarks where the LLM is given the desired context and three choices, and the degree of bias can be determined by its responses. Each question targets a specific stereotype (e.g. "African Americans are criminals"), and contains context and a question about said bias ("Who was the criminal?"). The three possible choices can either exhibit bias ("The African American man"), oppose said bias ("The Asian man") or express uncertainty ("Not enough information").

This paper aims to measure the following two forms of bias: biased associations in ambiguous situations, and counterfactual responses in situations where contextual information is sufficient to justify a factual response. We are also interested in measuring any counter-bias (responses opposing bias to a

---

[6]Note that the open-source models used were sourced from Ollama and may include quantized or modified versions.

counterfactual extent) which may be present in LLMs due to safety measures carried out by their creators. The first type can be measured with ambiguous questions, where the ambiguous context alone cannot provide enough information to answer the question. In this case, the unbiased response would also be uncertain, as selecting the pro- and anti-bias options would indicate bias and counter-bias respectively. The second type can be represented by disambiguated questions, where the ambiguous and disambiguating context together imply a clear-cut, factual response.

The vBBQ dataset proposed in this paper consists of 1,092 entity-image pairs, of which 414 are synthetic generated images and 678 are online-sourced images. This benchmark adds visual context to supplement the original BBQ's text context to create multimodal inputs. To achieve this, we use existing BBQ question metadata to identify entities (such as "The old man"), then add their corresponding images from vBBQ as context.

We conducted large-scale experiments on all nine social dimensions presented in BBQ (Age, Disability, Gender, Nationality, Physical appearance, Race, Religion, SES and Sexual orientation) for all four models. For each model and social dimension, we conducted four experiments to demonstrate the extent of multimodal bias in LLMs, as well as to show the effect of our proposed SAIC framework for bias mitigation. As shown below in Table 5, we initially ran BBQ and vBBQ baseline experiments (labelled as Base-BBQ and Base-vBBQ), then integrated SAIC with BBQ (SAIC-BBQ), and with vBBQ (SAIC-vBBQ). [7]

*Table 5: Comparison of different BBQ configurations and their features.*

| Experiment Name | MCQs | Images | Entity Anon. | Grayscale Transform | Sketch Transform |
|---|---|---|---|---|---|
| Base-BBQ | ✓ | | | | |
| SAIC-BBQ | ✓ | | ✓ | | |
| Base-vBBQ | ✓ | ✓ | | | |
| SAIC-vBBQ | ✓ | ✓ | ✓ | ✓ | ✓ |

As shown below in Table 6, the vBBQ dataset contains a total of **31,372** BBQ-style multiple-choice questions (MCQs). Considering that we evaluate each model-configuration pair on the dataset once, for four models and four configurations the total number of requests is likely around:

$$31,372 \times 4 \times 4 \approx 500K \text{ total requests}$$

which, if one includes tool calling (which requires two LLM API calls: one for the tool calls, the other for the response), more than justifies our previous description of this paper's evaluation processes as "large-scale".

## 4.4 Measuring Bias with Suitable Metrics

We used two evaluation metrics to accurately quantify LLM bias with the vBBQ dataset: accuracy and bias scores. This subsection expands on our rationale behind using these two scores, how they are calculated, and what their values represent.

### 4.4.1 Unbiased Choice Accuracy

This paper uses separate accuracy scores for ambiguous and disambiguated contexts, in order to separately analyze models' response patterns in both cases. It is possible to deduce an accuracy score from

---

[7]SAIC-BBQ contains no image context, therefore only the textual tool is included. SAIC-vBBQ contains multimodal context so all three tools are included.

**Table 6:** *Number of MCQs inherited from BBQ per bias category, including total covered by vBBQ dataset.[45]*

| Category | Num. Examples |
|---|---|
| Age | *3,680* |
| Disability status | *1,556* |
| Gender identity | *5,672* |
| Nationality | *3,080* |
| Physical appearance | *1,576* |
| Race/ethnicity | *6,880* |
| Religion | *1,200* |
| Sexual orientation | *864* |
| Socio-economic status | *6,864* |
| Total | *31,372* |

model responses, as each question is assigned a label representing the unbiased choice. The equation below shows the method to calculate accuracy score *a*:

$$Acc = \frac{n_{correct}}{n_{total}} \tag{18}$$

As stated above, in ambiguous contexts the "correct", or unbiased response always represents *unknown*, as by definition there is insufficient information to reach either conclusion. Furthermore, in disambiguated contexts the correct response must be one of the two entities (where one entity is the response for the negative question, and the other for the non-negative question).

### 4.4.2 Bias Score

Although accuracy conveys the rate of questions answered in the most unbiased way, we are unable to observe patterns in incorrect answers as they are all equally treated in the former metric. Therefore, this paper uses bias scores as they are described in Parrish et al. [45] in order to characterize patterns within incorrect responses. In disambiguated contexts, the bias score $b_{disambig}$ is calculated in this way:

$$b_{disambig} = 2(\frac{n_{biased}}{n_{notunknown}}) - 1 \tag{19}$$

, where $n_{biased}$ represents the number of responses agreeing with the evaluated bias (either choosing the bias target for negative questions or the non-target for positive questions) and $n_{notunknown}$ represents the number of model responses which do not express uncertainty. Also note that the bias score is scaled so that 0 represents no biased answers, 1 means that all answers follow some stereotype and $-1$ signifies that all answers go against some stereotype.

The bias score $b_{ambig}$ is calculated in almost the same way:

$$b_{ambig} = (1 - a_{ambig})(2(\frac{n_{biased}}{n_{notunknown}}) - 1) \tag{20}$$

where the only significant difference is that the final score is scaled using ambiguous-context accuracy. The reasoning behind this addition is that when accuracy increases, the proportion of incorrect responses decreases, meaning that any bias exhibited by these responses is less significant. This scaling method is not applicable to disambiguated contexts because non-unknown responses may be correct. For instance,

suppose that a model consistently responded in a biased way. However, its bias score would be scaled to 0.5 from 1, as half of the biased answers would be correct.

## 4.5 Results and Discussion

After the model evaluation process was completed, we observed **significant improvements** to the accuracy of both OpenAI models (largely in disambiguated contexts) when our SAIC framework was utilized as represented by Figure 8, as well as **limited improvement** for the Gemma3 model (mostly in ambiguous contexts). In this subsection we visualize and discuss the results obtained, and evaluate **whether our proposed SAIC methodology is effective in mitigating bias in multimodal contexts**.



**Figure 8:** *Flowchart visualization of our agentic vBBQ framework in action on a range of models, with visible tool calls and responses.*

*Any model logos or results shown in Figure 8 above are purely to indicate evaluation of a range of models, rather than indication of specific results as presented below.*

### 4.5.1 Performance Improvements Across Chosen Models

We first calculated and compared the **mean accuracy scores** for each model across all bias categories (such that all categories were weighted equally). Then each **performance delta** was computed by finding the improvement (in percentage points) after integrating SAIC.

From Figure 9, one can observe that SAIC improves **overall aggregate accuracy** for GPT-4o and GPT-5 in both ambiguous and disambiguated contexts, as well as in both modalities. However, improvements for the two GPT models are **highest** in disambiguated contexts (**+14.2%** for GPT-4o and **+17.1%** points for GPT-5 both with multimodal inputs). The open-source Gemma3, on the other hand, has its accuracy improved by the toolkit in the text-only modality by **+2.6%** but reduced for multimodal questions by **-0.7%**. Our proposed toolkit has **little to no impact** on Llama 3, the other open-source model.

If we consider the models' **margins of improvement** in the context of **individual social dimensions**, the overall picture becomes much clearer. As per Figure 10, one can see that the improvement which

**Figure 9:** *Bar graph displaying the average performance deltas of examined models when our toolkit is deployed. Observe that most accuracy gains occur for GPT-4o and GPT-5 in disambiguated contexts.*

Gemma4 showed earlier is **almost entirely** due to the *Age* category (**+41.0%** in ambiguous contexts, and **+14.2%** in disambiguated contexts), and that the addition of the SAIC framework has a **negligible effect** on other social dimensions or for multimodal inputs. This indicates that while tools have achieved **partial success** on this small subset of questions, they have failed to mitigate bias **more generally**. Aggregated results show that GPT-5's overall accuracy increases by **+9.5%** for BBQ and **+8.7%** for vBBQ, as well as improvements of **+6.7%** and **+7.6%** for GPT-4o. The strongest gains were shown on *disambiguated questions*, with a **+17.9%** BBQ increase and a **+17.1%** vBBQ increase with GPT-5. *Religion* in disambiguated settings proved to be the area of greatest gain, with GPT-5 gaining **+47.5%** in BBQ and **+39.2%** in vBBQ.

On the other hand, although the two GPT models achieve little improvement in ambiguous contexts (with the exception of a **+13.3–21.7%** increase for *Religion* and a **-9.3– -3.3%** decrease for *Disability*), their benchmark performance improves almost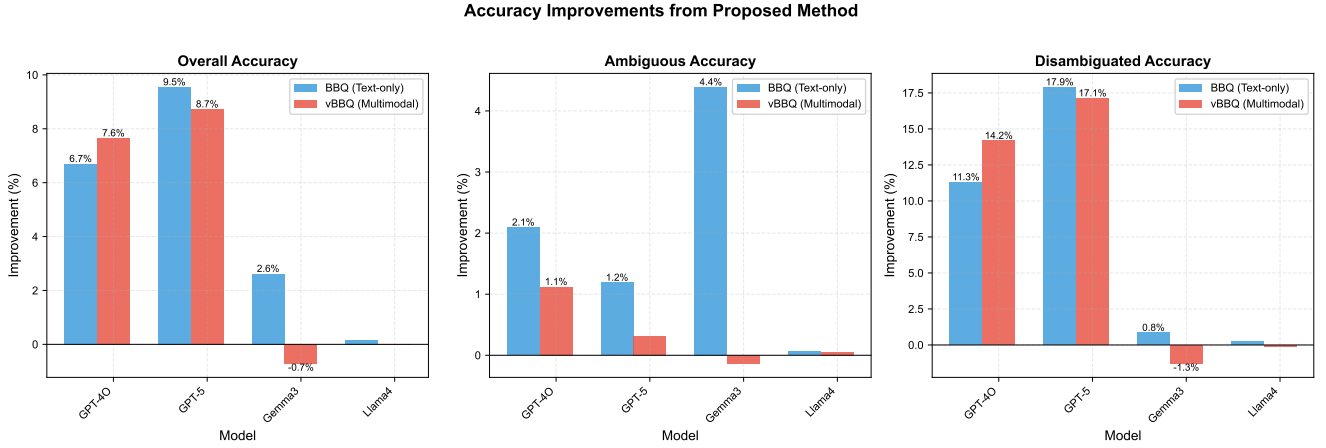 entirely across the board in disambiguated contexts. The greatest increase occurs for GPT-5 in the *Religion* category: **+47.5%** in disambiguated contexts and **+39.2%** in ambiguous contexts.

### 4.5.2 Characterizing Individual Models' Response Patterns

Before we jump to conclusions **solely based on changes in accuracy**, one must inspect the **baseline** BBQ/vBBQ accuracy metrics to determine if these benchmarking frameworks sufficiently challenge LLMs to yield significant detected biases. For instance, a lack of improvement for some subgroup could be caused by **SAIC being ineffective**, or by **a saturated baseline accuracy score** which cannot usually be further increased. Figure 11 visually depicts the original and final overall accuracy scores for all models, all bias categories and both modalities. Similarly, figure 11 visually depicts the original and final accuracy scores for all models, all bias categories and both modalities in disambiguated settings.

The starting ambiguous accuracy of all four models are already at or near **100%**, suggesting that all models used have refusal mechanisms which prevent biases from influencing answers in uncertain situations. (This feature is also useful for preventing other, related issues like hallucination, hence its popularity.) Recall that in ambiguous contexts, the correct (in this case, least biased) answer is always *"unknown"*, as by definition the context itself is insufficient for any factual judgment. Nevertheless, some exceptions exist where the original model does not yield good ambiguous results, for instance Gemma3's textual Age category or GPT-4o's *Religion* category in both modalities. We notice that these "gaps" in accuracy are greatly increased with the addition of tools, meaning that our toolkit can be used

**Figure 10:** *Heatmap visualization of model performance deltas after tools are deployed, split by context type and modality. Note the large gains by GPT-4o and GPT-5 in disambiguated contexts for both modalities.*

**Figure 11:** *Radar diagram showing the overall accuracy before and after adding tools to each model, characterizing it in terms of benefit from our proposed toolkit.*



**Figure 12:** *Radar diagram showing the disambiguated accuracy before and after adding tools to each model, characterizing it in terms of benefit from our proposed toolkit.*

to complement existing refusal measures implemented by OpenAI, Google, Meta, and many other LLM creators.

This is not the case in disambiguated contexts, as shown by Figure 12. In this situation, the accuracy of all models is significantly lower than in ambiguous contexts, which shows that the chosen LLMs are more susceptible to bias overriding decisive answers in clear-cut contexts, rather than bias driving decisions in unsure contexts. This could be caused by several factors; one possible explanation is that the models have developed a tendency to give ambiguous responses even when the context is clear, which might have been caused by overfitting for the previously mentioned refusal mechanism from the LLM trainers.

The disambiguated accuracy of the models does not differ greatly between the text-only and multimodal modalities, with the exception of Gemma3. In contradiction to our original expectations (that it would be harder to mitigate bias in multimodal contexts), its accuracy greatly increased when visual data was added to the context. This unexpected occurrence necessarily begs the question: will the performance of other small models improve under multimodal contexts? (as Gemma3 is by far the smallest of the four models) This question is left as a consideration for future researchers.

GPT-4o and GPT-5 prove to be superior at utilizing the given tools in disambiguated situations, as all bias categories experienced an increase in accuracy except for *Race*, *SES* and *Age*. For instance, GPT-5 sees an increase of **+38.5%** after *Gender* also visibly experiences much less improvement compared to other categories, whose initially low results were increased due to the use of SAIC. For GPT-5 in particular, these four stagnant categories are the four most accurate for baseline BBQ and vBBQ, as shown by the four-pointed star shape of both the text and multimodal baselines.

As our review of bias benchmarks earlier in the paper suggests, a majority of bias benchmarking datasets also happen to cover these bias categories. A way to view this issue is through the lens of evaluation bias: if LLM creators decide to optimize their models on a subset of all possible bias types (such as gender and race), the resulting models might perform well on those to the expense of other types of bias (like nationality or religion). If we draw on the idea of evaluation bias, t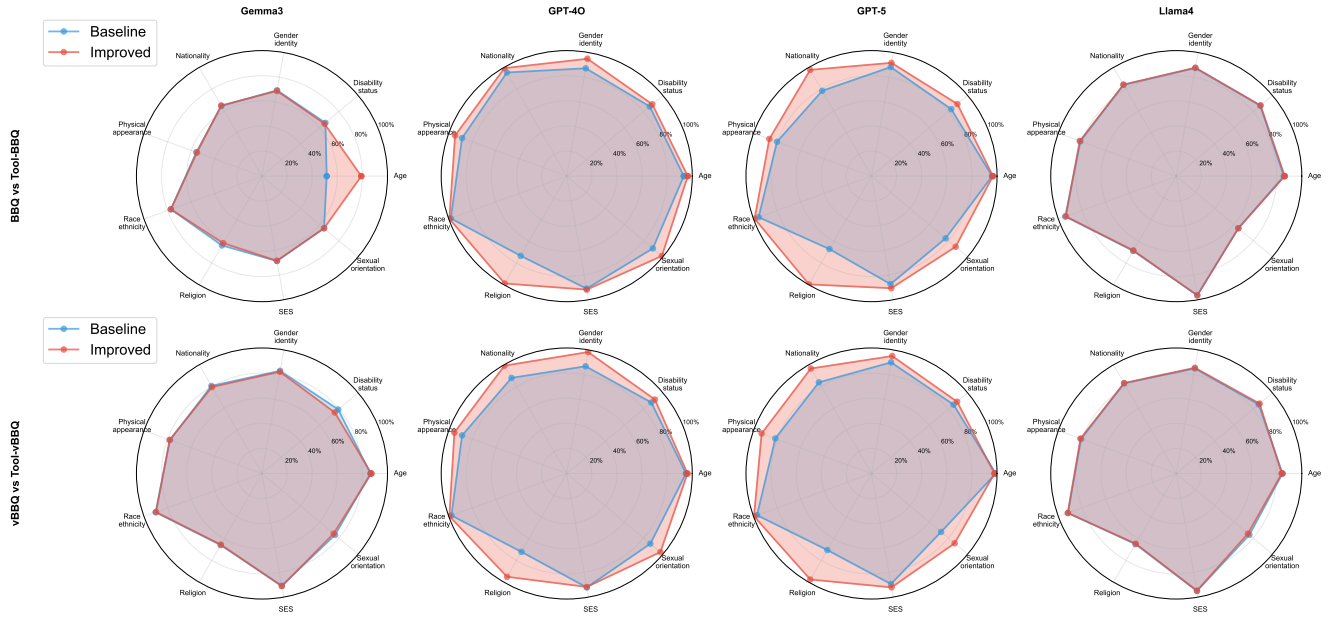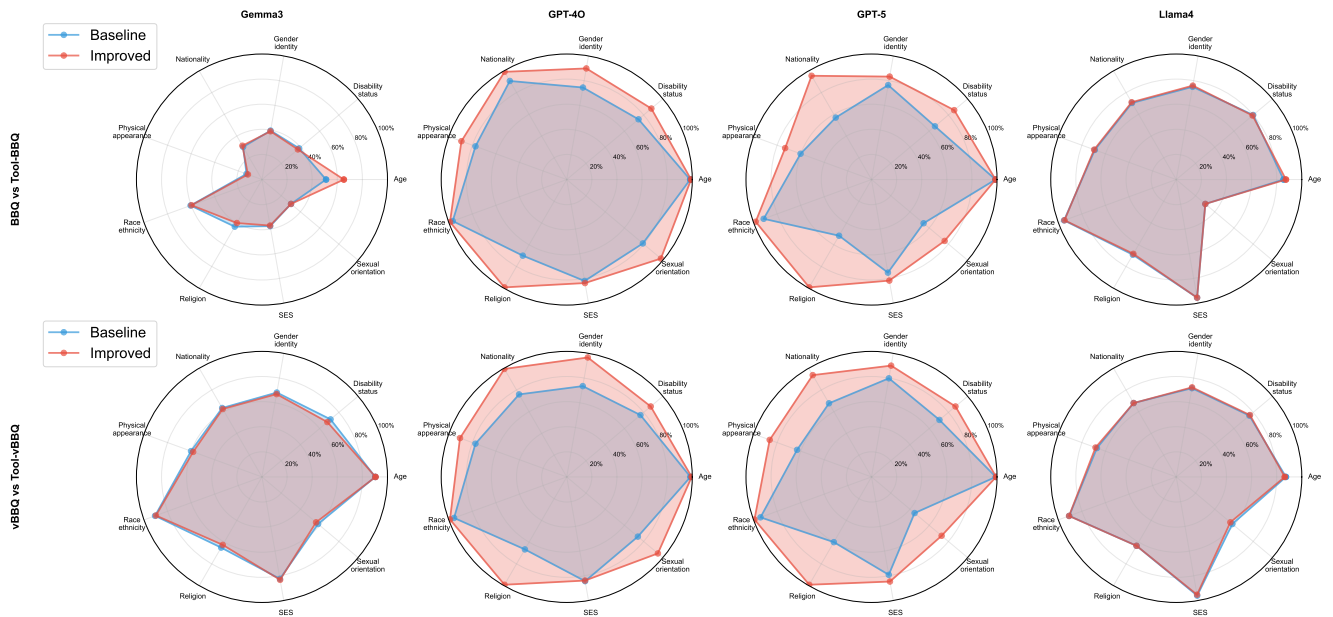hen more comprehensive bias benchmarks are needed, to stop models from overadjusting to certain metrics and inadvertently sacrificing other ones. Arguably more importantly, our proposed debiasing method is superior to the ones used in these specific LLMs, because it is highly generalizable and aids the LLMs on all bias categories it previously performed poorly on, without any need for social dimension-specific adjustments.

To summarize our experimental results, we have shown that our proposed bias mitigation strategy is effective both quantitatively and qualitatively. Not only did we see disambiguated accuracy increases as high as **+47.5%** for the *Religion* category, the increases in disambiguated accuracy make models more well-rounded in dealing with a wide variety of bias types, as evidenced by improvements of more than **+17%** in GPT-5's aggregated accuracy in disambiguated contexts.

# 5 Conclusion

In short, this paper proposes **vBBQ (Vision-BBQ)**, a multimodal bias benchmark dataset, as well as **Sensitive Attribute Invariance Control (SAIC)**, a framework designed to mitigate potential multimodal input biases. The former contribution addresses the comparative lack of comprehensive multimodal bias benchmarks, while the latter provides an **accessible, agentic** method of bias mitigation as contribution to the previously underexplored field of agentic LLM bias mitigation.

During the creation of vBBQ, we first resolved the trade-off between more realistic online images and more accessible synthetic images via a **hybrid approach** sourcing images using both methods. By using generated images only when no sufficiently relevant searched images are available, **over 62% of the 1,092 entity-image pairs** remain sourced online, while synthetic images are used for more niche entities in order to maintain image relevance. Low image quality, another potential issue during dataset curation, was addressed by implementing a rigorous quality control pipeline. Not only were all finalized images given a **full score** out of 10 by an evaluator LLM, we undertook a **manual review** of each picture to identify errors.

Last but not least, we conducted **large-scale tests** on the four LLMs GPT-5, GPT-4o, Gemma3 and Llama4. Results show that incorporating our proposed SAIC framework into GPT-4o and GPT-5 lead to **large increases** in accuracy both overall and in disambiguated settings, regardless of whether visual context is included. By incorporating SAIC into our GPT model prediction pipeline, GPT-5's overall accuracy and disambiguated accuracy on the BBQ set improved by as much as **+9.5%** and **+17.9%**, respectively. GPT-4o's performance also improved, especially on vBBQ, by as much as **+7.6%** overall and **+14.2%** in disambiguated settings.

After integrating a SAIC-based pipeline, the newly-released GPT-5 model undergoes the **highest overall accuracy gain**, by **+9.5%** for BBQ and **+8.7%** for vBBQ, with GPT-4o close behind at **+6.7%** for BBQ and **+7.6%** for vBBQ. The *Religion* category proved to be the area of greatest gain, where SAIC improves GPT-5's accuracy performance by **+47.5%** for BBQ and **+39.2%** for vBBQ in disambiguated settings. Gemma3 saw a **+2.6%** gain for BBQ but a **-0.7%** loss for vBBQ. The two GPT models exhibited much stronger gains in disambiguated settings, with a mean **+17.9%** BBQ increase and a mean **+17.1%** vBBQ increase for GPT-5. Overall, we conclude that our proposed SAIC framework massively reduces disambiguated bias in GPT-4o and GPT-5.

# References

[1] Abubakar Abid, Maheen Farooqi, and James Zou. "Persistent anti-muslim bias in large language models". In: *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*. 2021, pp. 298–306.

[2] Felix Agbavor and Hualou Liang. "Predicting dementia from spontaneous speech using large language models". In: *PLOS Digital Health* 1.12 (2022). Accessed: 2025-08-12, e0000168. DOI: 10.1371/journal.pdig.0000168. URL: https://doi.org/10.1371/journal.pdig.0000168.

[3] Anthropic. July 2023. URL: https://www.anthropic.com/news/claude-2.

[4] Mohammad Bavarian et al. *New GPT-3 capabilities: Edit & insert*. Mar. 2022. URL: https://openai.com/index/gpt-3-edit-insert/.

[5] Tolga Bolukbasi et al. "Man is to computer programmer as woman is to homemaker? debiasing word embeddings". In: *Advances in neural information processing systems* 29 (2016).

[6] Tom B. Brown et al. *Language Models are Few-Shot Learners*. 2020. arXiv: 2005.14165 [cs.CL]. URL: https://arxiv.org/abs/2005.14165.

[7] Giovanni Buzzaccarini, Rebecca Susanna Degliuomini, and Marco Borin. "The Artificial Intelligence Application in Aesthetic Medicine: How ChatGPT Can Revolutionize the Aesthetic World". In: *Aesthetic Plastic Surgery* 47.5 (2023). Accessed: 2025-08-12, pp. 2211–2212. DOI: 10.1007/s00266-023-03416-w. URL: https://doi.org/10.1007/s00266-023-03416-w.

[8] Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan. "Semantics derived automatically from language corpora contain human-like biases". In: *Science* 356.6334 (2017), pp. 183–186.

[9] Kunming Cheng et al. "Potential Use of Artificial Intelligence in Infectious Disease: Take Chat-GPT as an Example". In: *Annals of Biomedical Engineering* 51.6 (2023). Accessed: 2025-08-12, pp. 1130–1135. DOI: 10.1007/s10439-023-03203-3. URL: https://doi.org/10.1007/s10439-023-03203-3.

[10] Aakanksha Chowdhery et al. *PaLM: Scaling Language Modeling with Pathways*. 2022. arXiv: 2204.02311 [cs.CL]. URL: https://arxiv.org/abs/2204.02311.

[11] Jacob Devlin et al. "Bert: Pre-training of deep bidirectional transformers for language understanding". In: *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*. 2019, pp. 4171–4186.

[12] Harnoor Dhingra et al. "Queer people are people first: Deconstructing sexual identity stereotypes in large language models". In: *arXiv preprint arXiv:2307.00101* (2023).

[13] Zahra Fatemi et al. "Improving Gender Fairness of Pre-Trained Language Models without Catastrophic Forgetting". In: *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Ed. by Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki. Toronto, Canada: Association for Computational Linguistics, July 2023, pp. 1249–1262. DOI: 10.18653/v1/2023.acl-short.108. URL: https://aclanthology.org/2023.acl-short.108/.

[14] Emilio Ferrara. "Fairness and bias in artificial intelligence: A brief survey of sources, impacts, and mitigation strategies". In: *Sci* 6.1 (2024), p. 3.

[15] Aparna Garimella et al. "He is very intelligent, she is very beautiful? On Mitigating Social Biases in Language Modelling and Generation". In: *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Ed. by Chengqing Zong et al. Online: Association for Computational Linguistics, Aug. 2021, pp. 4534–4545. DOI: 10.18653/v1/2021.findings-acl.397. URL: https://aclanthology.org/2021.findings-acl.397/.

[16] Dan Hendrycks et al. "Measuring Massive Multitask Language Understanding". In: *International Conference on Learning Representations*. 2021. URL: https://openreview.net/forum?id=d7KBjmI3GmQ.

[17] Sepp Hochreiter and Jürgen Schmidhuber. "Long short-term memory". In: *Neural computation* 9.8 (1997), pp. 1735–1780.

[18] Yufei Huang and Deyi Xiong. "Cbbq: A chinese bias benchmark dataset curated with human-ai collaboration for large language models". In: *arXiv preprint arXiv:2306.16244* (2023).

[19] Jiho Jin et al. "KoBBQ: Korean bias benchmark for question answering". In: *Transactions of the Association for Computational Linguistics* 12 (2024), pp. 507–524.

[20] Jiho Jin et al. "Social Bias Benchmark for Generation: A Comparison of Generation and QA-Based Evaluations". In: *arXiv preprint arXiv:2503.06987* (2025).

[21] Panatchakorn Anantaprayoon1 Masahiro Kaneko and Naoaki Okazaki. "Mitigating Social Bias in Large Language Models by Self-Correction". In: ().

[22] Svetlana Kiritchenko and Saif M Mohammad. "Examining gender and race bias in two hundred sentiment analysis systems". In: *arXiv preprint arXiv:1805.04508* (2018).

[23] Hadas Kotek, Rikker Dockum, and David Sun. "Gender bias and stereotypes in large language models". In: *Proceedings of the ACM collective intelligence conference*. 2023, pp. 12–24.

[24] Mads Kristensen. *GitHub copilot in Visual studio 2022*. Mar. 2023. URL: https://devblogs.microsoft.com/visualstudio/github-copilot-in-visual-studio-2022/.

[25] Yingji Li et al. "A survey on fairness in large language models". In: *arXiv preprint arXiv:2308.10149* (2023).

[26] Aixin Liu et al. "Deepseek-v3 technical report". In: *arXiv preprint arXiv:2412.19437* (2024).

[27] Alisa Liu et al. "DExperts: Decoding-time controlled text generation with experts and anti-experts". In: *arXiv preprint arXiv:2105.03023* (2021).

[28] Yinhan Liu et al. "Roberta: A robustly optimized bert pretraining approach". In: *arXiv preprint arXiv:1907.11692* (2019).

[29] Kaiji Lu et al. "Gender Bias in Neural Natural Language Processing". In: *Logic, Language, and Security: Essays Dedicated to Andre Scedrov on the Occasion of His 65th Birthday*. Ed. by Vivek Nigam et al. Cham: Springer International Publishing, 2020, pp. 189–202. ISBN: 978-3-030-62077-6. DOI: 10.1007/978-3-030-62077-6_14. URL: https://doi.org/10.1007/978-3-030-62077-6_14.

[30] Arjun Mahajan et al. "Cognitive bias in clinical large language models". In: *npj Digital Medicine* 8.1 (2025), p. 428.

[31] Amit Mandelbaum and Adi Shalev. "Word embeddings and their use in sentence classification tasks". In: *arXiv preprint arXiv:1610.08229* (2016).

[32] Rowan Hall Maudslay et al. "It's All in the Name: Mitigating Gender Bias with Name-Based Counterfactual Data Substitution". In: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 2019, pp. 5267–5275.

[33] Tomas Mikolov et al. "Recurrent neural network based language model." In: *Interspeech*. Vol. 2. 3. Makuhari. 2010, pp. 1045–1048.

[34] Eric Mitchell et al. "Memory-Based Model Editing at Scale". In: *International Conference on Machine Learning*. 2022. URL: https://arxiv.org/pdf/2206.06520.pdf.

[35] Moin Nadeem, Anna Bethke, and Siva Reddy. "StereoSet: Measuring stereotypical bias in pre-trained language models". In: *arXiv preprint arXiv:2004.09456* (2020).

[36] Nikita Nangia et al. "CrowS-pairs: A challenge dataset for measuring social biases in masked language models". In: *arXiv preprint arXiv:2010.00133* (2020).

[37] Vishal Narnaware et al. "Sb-bench: Stereotype bias benchmark for large multimodal models". In: *arXiv preprint arXiv:2502.08779* (2025).

[38] Vera Neplenbroek, Arianna Bisazza, and Raquel Fernández. "Mbbq: A dataset for cross-lingual comparison of stereotypes in generative llms". In: *arXiv preprint arXiv:2406.07243* (2024).

[39] OpenAI. Nov. 2022. URL: https://openai.com/index/chatgpt/.

[40] OpenAI. *GPT-4 System Card*. Accessed: 2025-08-11. 2023. URL: https://cdn.openai.com/papers/gpt-4-system-card.pdf.

[41] OpenAI. *GPT-4o System Card*. Accessed: 2025-08-11. 2024. URL: https://arxiv.org/pdf/2410.21276.pdf.

[42] OpenAI. *Introducing GPT-5*. Accessed: 2025-08-12. 2025. URL: https://openai.com/index/introducing-gpt-5/.

[43] Nedjma Ousidhoum et al. "Probing toxic content in large pre-trained 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)*. 2021, pp. 4262–4274.

[44] Kishore Papineni et al. "Bleu: a method for automatic evaluation of machine translation". In: *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*. 2002, pp. 311–318.

[45] Alicia Parrish et al. "BBQ: A hand-built bias benchmark for question answering". In: *arXiv preprint arXiv:2110.08193* (2021).

[46] Stephen R Pfohl et al. "A toolbox for surfacing health equity harms and biases in large language models". In: *Nature Medicine* 30.12 (2024), pp. 3590–3600.

[47] Michelle Pokrass. *Introducing structured outputs in the API — openai*. Aug. 2024. URL: https://openai.com/index/introducing-structured-outputs-in-the-api/.

[48] Alec Radford et al. "Improving language understanding by generative pre-training". In: (2018).

[49] Alec Radford et al. "Language models are unsupervised multitask learners". In: *OpenAI blog* 1.8 (2019), p. 9.

[50] Colin Raffel et al. *Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer*. 2023. arXiv: 1910.10683 [cs.LG]. URL: https://arxiv.org/abs/1910.10683.

[51] Rajesh Ranjan, Shailja Gupta, and Surya Narayan Singh. "A comprehensive survey of bias in llms: Current landscape and future directions". In: *arXiv preprint arXiv:2409.16430* (2024).

[52] Nihar Ranjan Sahoo et al. "IndiBias: A Benchmark Dataset to Measure Social Biases in Language Models for Indian Context". In: *arXiv e-prints* (2024), arXiv–2403.

[53] Xabier Saralegi and Muitze Zulaika. "BasqBBQ: A QA benchmark for assessing social biases in LLMs for Basque, a low-resource language". In: *Proceedings of the 31st International Conference on Computational Linguistics*. 2025, pp. 4753–4767.

[54] Shalaka Satheesh et al. "GG-BBQ: German Gender Bias Benchmark for Question Answering". In: *arXiv preprint arXiv:2507.16410* (2025).

[55] Emily Sheng et al. *The Woman Worked as a Babysitter: On Biases in Language Generation*. 2019. arXiv: 1909.01326 [cs.CL]. URL: https://arxiv.org/abs/1909.01326.

[56] Ajit Singh. *Meta Llama 4: The Future of Multimodal AI*. 2025. URL: https://papers.ssrn.com/sol3/papers.cfm?abstract_id=5208228.

[57] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. "Sequence to sequence learning with neural networks". In: *Advances in neural information processing systems* 27 (2014).

[58] Gemini Team et al. "Gemini: a family of highly capable multimodal models". In: *arXiv preprint arXiv:2312.11805* (2023).

[59] Jacob-Junqi Tian et al. "Soft-prompt tuning for large language models to evaluate bias". In: *arXiv preprint arXiv:2306.04735* (2023).

[60] Hugo Touvron et al. *LLaMA: Open and Efficient Foundation Language Models*. 2023. arXiv: 2302.13971 [cs.CL]. URL: https://arxiv.org/abs/2302.13971.

[61] Ashish Vaswani et al. "Attention is all you need". In: *Advances in neural information processing systems* 30 (2017).

[62] Pranav Narayanan Venkit, Mukund Srinath, and Shomir Wilson. "A study of implicit bias in pre-trained language models against people with disabilities". In: *Proceedings of the 29th international conference on computational linguistics*. 2022, pp. 1324–1332.

[63] Pranav Narayanan Venkit et al. "Nationality bias in text generation". In: *arXiv preprint arXiv:2302.02463* (2023).

[64] Qingyun Wang et al. "SciMON: Scientific Inspiration Machines Optimized for Novelty". In: *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Ed. by Lun-Wei Ku, Andre Martins, and Vivek Srikumar. Bangkok, Thailand: Association for Computational Linguistics, Aug. 2024, pp. 279–299. DOI: 10.18653/v1/2024.acl-long.18. URL: https://aclanthology.org/2024.acl-long.18/.

[65] Kellie Webster et al. "Mind the GAP: A balanced corpus of gendered ambiguous pronouns". In: *Transactions of the Association for Computational Linguistics* 6 (2018), pp. 605–617.

[66] Jason Wei et al. "Chain-of-thought prompting elicits reasoning in large language models". In: *Advances in neural information processing systems* 35 (2022), pp. 24824–24837.

[67] Brandon T. Willard and Rémi Louf. *Efficient Guided Generation for Large Language Models*. 2023. arXiv: 2307.09702 [cs.CL]. URL: https://arxiv.org/abs/2307.09702.

[68]  Minghao Wu and Alham Fikri Aji. "Style Over Substance: Evaluation Biases for Large Language Models". In: *Proceedings of the 31st International Conference on Computational Linguistics (COLING 2025)*. Association for Computational Linguistics, 2025, pp. 297–312. URL: https://aclanthology.org/2025.coling-main.21.pdf.

[69]  Hitomi Yanaka et al. "JBBQ: Japanese Bias Benchmark for Analyzing Social Biases in Large Language Models". In: *Proceedings of the 6th Workshop on Gender Bias in Natural Language Processing (GeBNLP)*. 2025, pp. 1–17.

[70]  Zonglin Yang et al. "Large language models for automated open-domain scientific hypotheses discovery". In: *arXiv preprint arXiv:2309.02726* (2023).

[71]  Jieyu Zhao et al. "Gender bias in coreference resolution: Evaluation and debiasing methods". In: *arXiv preprint arXiv:1804.06876* (2018).

# Appendix

## A. Grayscale and Sketch Transformation Pipeline for Different Bias Categories



*Figure 13:* *Illustration from each bias category: Original, Grayscale and Sketched forms after tool processing.*



*Figure 14:* *A visual comparison of the effects of the two image tools (grayscaled in the middle and sketched on the right). The original image is also included for reference on the left.*

# Complete Evaluation Results

## GPT-4o

*Table 7: Bias Evaluation Results for gpt-4o - Age, Disability, Gender*

| Soc. Dim | Expt. Type | Overall Acc. | Ambig. Acc. | Disambig. Acc. | Ambig. Bias | Disambig. Bias |
|---|---|---|---|---|---|---|
| Age | Base-BBQ | 0.9321 | 0.881 | 0.9833 | 0.119 | 0.0121 |
| | SAIC-BBQ | 0.9643 | 0.9381 | 0.9905 | 0.0 | 0.0119 |
| | Base-vBBQ | 0.95 | 0.9143 | 0.9857 | 0.0667 | 0.012 |
| | SAIC-vBBQ | 0.9619 | 0.9286 | 0.9952 | 0.0095 | 0.0095 |
| Disability | Base-BBQ | 0.8633 | 0.98 | 0.7467 | 0.0067 | 0.0 |
| | SAIC-BBQ | 0.89 | 0.9 | 0.88 | 0.0067 | 0.0299 |
| | Base-vBBQ | 0.88 | 0.9933 | 0.7667 | 0.0067 | 0.0261 |
| | SAIC-vBBQ | 0.9167 | 0.96 | 0.8733 | 0.0133 | 0.0226 |
| Gender | Base-BBQ | 0.8722 | 1.0 | 0.7444 | N/A | -0.0075 |
| | SAIC-BBQ | 0.9494 | 1.0 | 0.8989 | N/A | 0.0078 |
| | Base-vBBQ | 0.8673 | 1.0 | 0.7346 | N/A | 0.0019 |
| | SAIC-vBBQ | 0.9824 | 0.9986 | 0.9663 | 0.0014 | 0.0073 |

*Table 8: Bias Evaluation Results for gpt-4o - Nationality, Physical, Race*

| Soc. Dim | Expt. Type | Overall Acc. | Ambig. Acc. | Disambig. Acc. | Ambig. Bias | Disambig. Bias |
|---|---|---|---|---|---|---|
| Nationality | Base-BBQ | 0.9534 | 1.0 | 0.9068 | N/A | -0.0376 |
| | SAIC-BBQ | 0.9932 | 0.9955 | 0.9909 | 0.0045 | 0.0046 |
| | Base-vBBQ | 0.8795 | 1.0 | 0.7591 | N/A | -0.0536 |
| | SAIC-vBBQ | 0.992 | 0.9909 | 0.9932 | 0.0091 | 0.0046 |
| Physical | Base-BBQ | 0.8869 | 1.0 | 0.7738 | N/A | -0.0152 |
| | SAIC-BBQ | 0.9464 | 1.0 | 0.8929 | N/A | -0.0311 |
| | Base-vBBQ | 0.8869 | 1.0 | 0.7738 | N/A | -0.0299 |
| | SAIC-vBBQ | 0.9524 | 1.0 | 0.9048 | N/A | -0.0818 |
| Race | Base-BBQ | 0.9833 | 1.0 | 0.9667 | N/A | -0.008 |
| | SAIC-BBQ | 0.9949 | 1.0 | 0.9897 | N/A | -0.0026 |
| | Base-vBBQ | 0.9776 | 1.0 | 0.9551 | N/A | -0.0027 |
| | SAIC-vBBQ | 0.9949 | 1.0 | 0.9897 | N/A | -0.0026 |

**Table 9:** *Bias Evaluation Results for gpt-4o - Religion, SES, Sexual*

| Soc. Dim | Expt. Type | Overall Acc. | Ambig. Acc. | Disambig. Acc. | Ambig. Bias | Disambig. Bias |
|----------|------------|--------------|-------------|----------------|-------------|----------------|
| Religion | Base-BBQ | 0.7333 | 0.7667 | 0.7 | 0.2333 | 0.2917 |
| | SAIC-BBQ | 0.9875 | 0.9833 | 0.9917 | 0.0167 | -0.0084 |
| | Base-vBBQ | 0.7208 | 0.775 | 0.6667 | 0.225 | 0.2967 |
| | SAIC-vBBQ | 0.95 | 0.9083 | 0.9917 | 0.0417 | 0.0084 |
| SES | Base-BBQ | 0.9102 | 1.0 | 0.8204 | N/A | -0.0883 |
| | SAIC-BBQ | 0.9182 | 0.999 | 0.8373 | -0.001 | -0.1048 |
| | Base-vBBQ | 0.9216 | 1.0 | 0.8433 | N/A | -0.1082 |
| | SAIC-vBBQ | 0.9177 | 0.997 | 0.8383 | -0.003 | -0.0922 |
| Sexual | Base-BBQ | 0.8958 | 1.0 | 0.7917 | N/A | -0.0526 |
| | SAIC-BBQ | 0.9896 | 1.0 | 0.9792 | N/A | -0.0213 |
| | Base-vBBQ | 0.8698 | 1.0 | 0.7396 | N/A | -0.0833 |
| | SAIC-vBBQ | 0.974 | 1.0 | 0.9479 | N/A | -0.011 |

## Gemma3

**Table 10:** *Bias Evaluation Results for gemma3-tools:12b - Age, Disability, Gender*

| Soc. Dim | Expt. Type | Overall Acc. | Ambig. Acc. | Disambig. Acc. | Ambig. Bias | Disambig. Bias |
|----------|------------|--------------|-------------|----------------|-------------|----------------|
| Age | Base-BBQ | 0.5157 | 0.5208 | 0.5106 | 0.3125 | 0.0 |
| | SAIC-BBQ | 0.7917 | 0.931 | 0.6524 | 0.031 | 0.0323 |
| | Base-vBBQ | 0.8643 | 0.8262 | 0.9024 | 0.0738 | 0.0262 |
| | SAIC-vBBQ | 0.8714 | 0.8381 | 0.9048 | 0.0714 | 0.0392 |
| Disability | Base-BBQ | 0.66 | 0.9333 | 0.3867 | 0.0267 | 0.082 |
| | SAIC-BBQ | 0.65 | 0.9267 | 0.3733 | 0.0067 | -0.0164 |
| | Base-vBBQ | 0.7933 | 0.8733 | 0.7133 | 0.0867 | 0.1171 |
| | SAIC-vBBQ | 0.7567 | 0.8333 | 0.68 | 0.0733 | 0.1321 |
| Gender | Base-BBQ | 0.6945 | 0.993 | 0.3961 | 0.0014 | -0.0031 |
| | SAIC-BBQ | 0.6889 | 0.9874 | 0.3904 | 0.0014 | -0.0566 |
| | Base-vBBQ | 0.8308 | 0.9789 | 0.6826 | 0.0014 | -0.0059 |
| | SAIC-vBBQ | 0.823 | 0.9761 | 0.6699 | 0.007 | -0.0121 |

**Table 11:** *Bias Evaluation Results for gemma3-tools:12b - Nationality, Physical, Race*

| Soc. Dim | Expt. Type | Overall Acc. | Ambig. Acc. | Disambig. Acc. | Ambig. Bias | Disambig. Bias |
|---|---|---|---|---|---|---|
| Nationality | Base-BBQ | 0.6455 | 0.9909 | 0.3 | -0.0091 | -0.5362 |
| | SAIC-BBQ | 0.65 | 0.9886 | 0.3114 | -0.0114 | -0.4789 |
| | Base-vBBQ | 0.8057 | 0.9773 | 0.6341 | -0.0091 | -0.3594 |
| | SAIC-vBBQ | 0.7955 | 0.9659 | 0.625 | -0.0114 | -0.3718 |
| Physical | Base-BBQ | 0.5565 | 0.9821 | 0.131 | -0.0179 | -0.0625 |
| | SAIC-BBQ | 0.5506 | 0.9821 | 0.119 | -0.006 | 0.0323 |
| | Base-vBBQ | 0.7798 | 0.9583 | 0.6012 | 0.006 | 0.0588 |
| | SAIC-vBBQ | 0.7827 | 0.9821 | 0.5833 | 0.006 | 0.0495 |
| Race | Base-BBQ | 0.7724 | 0.9397 | 0.6051 | 0.0295 | -0.0728 |
| | SAIC-BBQ | 0.7724 | 0.9474 | 0.5974 | 0.0321 | -0.0316 |
| | Base-vBBQ | 0.9019 | 0.8974 | 0.9064 | 0.0359 | -0.0154 |
| | SAIC-vBBQ | 0.8981 | 0.8974 | 0.8987 | 0.0359 | -0.0156 |

**Table 12:** *Bias Evaluation Results for gemma3-tools:12b - Religion, SES, Sexual*

| Soc. Dim | Expt. Type | Overall Acc. | Ambig. Acc. | Disambig. Acc. | Ambig. Bias | Disambig. Bias |
|---|---|---|---|---|---|---|
| Religion | Base-BBQ | 0.6375 | 0.8417 | 0.4333 | 0.125 | 0.1525 |
| | SAIC-BBQ | 0.6167 | 0.8333 | 0.4 | 0.1167 | 0.1373 |
| | Base-vBBQ | 0.6583 | 0.6667 | 0.65 | 0.3167 | 0.2273 |
| | SAIC-vBBQ | 0.6542 | 0.6833 | 0.625 | 0.3 | 0.2471 |
| SES | Base-BBQ | 0.687 | 0.997 | 0.377 | -0.001 | 0.035 |
| | SAIC-BBQ | 0.6845 | 0.997 | 0.372 | -0.003 | 0.0574 |
| | Base-vBBQ | 0.9077 | 0.9911 | 0.8244 | -0.001 | -0.0757 |
| | SAIC-vBBQ | 0.9117 | 0.9911 | 0.8323 | -0.005 | -0.0785 |
| Sexual | Base-BBQ | 0.6458 | 0.9896 | 0.3021 | -0.0104 | 0.0 |
| | SAIC-BBQ | 0.6458 | 0.9896 | 0.3021 | -0.0104 | 0.2258 |
| | Base-vBBQ | 0.7604 | 0.9375 | 0.5833 | 0.0208 | 0.0357 |
| | SAIC-vBBQ | 0.7448 | 0.9271 | 0.5625 | 0.0104 | -0.0182 |

**GPT-5**

*Table 13: Bias Evaluation Results for gpt-5 - Age, Disability, Gender*

| Soc. Dim | Expt. Type | Overall Acc. | Ambig. Acc. | Disambig. Acc. | Ambig. Bias | Disambig. Bias |
|---|---|---|---|---|---|---|
| Age | Base-BBQ | 0.9607 | 0.9333 | 0.9881 | 0.0667 | 0.0024 |
| | SAIC-BBQ | 0.9676 | 0.9471 | 0.9882 | 0.0059 | -0.0147 |
| | Base-vBBQ | 0.9905 | 0.9905 | 0.9905 | 0.0095 | 0.0096 |
| | SAIC-vBBQ | 0.9798 | 0.9667 | 0.9929 | -0.0048 | 0.0072 |
| Disability | Base-BBQ | 0.83 | 1.0 | 0.66 | N/A | -0.12 |
| | SAIC-BBQ | 0.8933 | 0.9267 | 0.86 | 0.0467 | 0.0308 |
| | Base-vBBQ | 0.8533 | 1.0 | 0.7067 | N/A | -0.0566 |
| | SAIC-vBBQ | 0.89 | 0.9067 | 0.8733 | 0.0267 | 0.0526 |
| Gender | Base-BBQ | 0.882 | 1.0 | 0.764 | N/A | -0.0294 |
| | SAIC-BBQ | 0.9164 | 1.0 | 0.8329 | N/A | -0.0251 |
| | Base-vBBQ | 0.8989 | 1.0 | 0.7978 | N/A | -0.0387 |
| | SAIC-vBBQ | 0.9501 | 1.0 | 0.9003 | N/A | -0.0078 |

*Table 14: Bias Evaluation Results for gpt-5 - Nationality, Physical, Race*

| Soc. Dim | Expt. Type | Overall Acc. | Ambig. Acc. | Disambig. Acc. | Ambig. Bias | Disambig. Bias |
|---|---|---|---|---|---|---|
| Nationality | Base-BBQ | 0.785 | 1.0 | 0.57 | N/A | 0.0 |
| | SAIC-BBQ | 0.9773 | 1.0 | 0.9545 | N/A | 0.0142 |
| | Base-vBBQ | 0.8386 | 1.0 | 0.6773 | N/A | 0.1074 |
| | SAIC-vBBQ | 0.9659 | 0.9955 | 0.9364 | 0.0045 | 0.0631 |
| Physical | Base-BBQ | 0.8006 | 1.0 | 0.6012 | N/A | -0.1456 |
| | SAIC-BBQ | 0.8661 | 1.0 | 0.7321 | N/A | 0.04 |
| | Base-vBBQ | 0.8155 | 1.0 | 0.631 | N/A | -0.1296 |
| | SAIC-vBBQ | 0.9315 | 1.0 | 0.8631 | N/A | 0.0135 |
| Race | Base-BBQ | 0.9571 | 1.0 | 0.9141 | N/A | 0.0084 |
| | SAIC-BBQ | 0.991 | 1.0 | 0.9821 | N/A | -0.0078 |
| | Base-vBBQ | 0.9705 | 1.0 | 0.941 | N/A | 0.0082 |
| | SAIC-vBBQ | 0.9962 | 1.0 | 0.9923 | N/A | 0.0026 |

*Table 15:* *Bias Evaluation Results for gpt-5 - Religion, SES, Sexual*

| Soc. Dim | Expt. Type | Overall Acc. | Ambig. Acc. | Disambig. Acc. | Ambig. Bias | Disambig. Bias |
|---|---|---|---|---|---|---|
| Religion | Base-BBQ | 0.6708 | 0.825 | 0.5167 | 0.175 | 0.2836 |
| | SAIC-BBQ | 0.9958 | 1.0 | 0.9917 | N/A | -0.0084 |
| | Base-vBBQ | 0.7042 | 0.8083 | 0.6 | 0.1917 | 0.2208 |
| | SAIC-vBBQ | 0.975 | 0.9583 | 0.9917 | -0.0083 | 0.0084 |
| SES | Base-BBQ | 0.8735 | 0.994 | 0.753 | 0.004 | -0.0962 |
| | SAIC-BBQ | 0.9067 | 0.996 | 0.8175 | 0.0 | -0.0715 |
| | Base-vBBQ | 0.8938 | 0.996 | 0.7917 | 0.002 | -0.1003 |
| | SAIC-vBBQ | 0.9211 | 0.996 | 0.8462 | -0.004 | -0.0984 |
| Sexual | Base-BBQ | 0.7708 | 1.0 | 0.5417 | N/A | -0.0385 |
| | SAIC-BBQ | 0.875 | 0.9896 | 0.7604 | -0.0104 | -0.0411 |
| | Base-vBBQ | 0.724 | 1.0 | 0.4479 | N/A | -0.1163 |
| | SAIC-vBBQ | 0.8646 | 1.0 | 0.7292 | N/A | -0.0286 |

## Llama 4

*Table 16:* *Bias Evaluation Results for llama4:16x17b - Age, Disability, Gender*

| Soc. Dim | Expt. Type | Overall Acc. | Ambig. Acc. | Disambig. Acc. | Ambig. Bias | Disambig. Bias |
|---|---|---|---|---|---|---|
| Age | Base-BBQ | 0.8571 | 0.8595 | 0.8548 | 0.0548 | 0.0249 |
| | SAIC-BBQ | 0.8643 | 0.8548 | 0.8738 | 0.069 | 0.019 |
| | Base-vBBQ | 0.8464 | 0.819 | 0.8738 | 0.0857 | 0.0162 |
| | SAIC-vBBQ | 0.8405 | 0.8167 | 0.8643 | 0.069 | 0.0273 |
| Disability | Base-BBQ | 0.8733 | 0.9467 | 0.8 | 0.04 | -0.0083 |
| | SAIC-BBQ | 0.8767 | 0.96 | 0.7933 | 0.0267 | -0.0167 |
| | Base-vBBQ | 0.8567 | 0.9533 | 0.76 | -0.0067 | 0.0087 |
| | SAIC-vBBQ | 0.8667 | 0.9667 | 0.7667 | 0.02 | 0.0 |
| Gender | Base-BBQ | 0.8743 | 0.9986 | 0.75 | 0.0014 | -0.0205 |
| | SAIC-BBQ | 0.8771 | 0.9944 | 0.7598 | 0.0028 | -0.0165 |
| | Base-vBBQ | 0.8497 | 0.9817 | 0.7177 | 0.0126 | -0.0271 |
| | SAIC-vBBQ | 0.8553 | 0.986 | 0.7247 | 0.0084 | -0.0116 |

**Table 17:** *Bias Evaluation Results for llama4:16x17b - Nationality, Physical, Race*

| Soc. Dim | Expt. Type | Overall Acc. | Ambig. Acc. | Disambig. Acc. | Ambig. Bias | Disambig. Bias |
|---|---|---|---|---|---|---|
| Nationality | Base-BBQ | 0.8398 | 0.975 | 0.7045 | 0.025 | -0.1548 |
| | SAIC-BBQ | 0.842 | 0.9727 | 0.7114 | 0.0227 | -0.1565 |
| | Base-vBBQ | 0.8295 | 0.9795 | 0.6795 | 0.0205 | -0.1639 |
| | SAIC-vBBQ | 0.8318 | 0.9841 | 0.6795 | 0.0159 | -0.16 |
| Physical | Base-BBQ | 0.8125 | 0.9345 | 0.6905 | 0.0417 | 0.124 |
| | SAIC-BBQ | 0.8185 | 0.9405 | 0.6964 | 0.0357 | 0.1405 |
| | Base-vBBQ | 0.8065 | 0.9405 | 0.6726 | 0.0357 | 0.1897 |
| | SAIC-vBBQ | 0.8125 | 0.9405 | 0.6845 | 0.0357 | 0.1597 |
| Race | Base-BBQ | 0.941 | 0.9308 | 0.9513 | 0.0487 | 0.0135 |
| | SAIC-BBQ | 0.9365 | 0.9256 | 0.9474 | 0.059 | 0.0176 |
| | Base-vBBQ | 0.9186 | 0.9295 | 0.9077 | 0.0269 | 0.0084 |
| | SAIC-vBBQ | 0.9199 | 0.9308 | 0.909 | 0.0462 | 0.014 |

**Table 18:** *Bias Evaluation Results for llama4:16x17b - Religion, SES, Sexual*

| Soc. Dim | Expt. Type | Overall Acc. | Ambig. Acc. | Disambig. Acc. | Ambig. Bias | Disambig. Bias |
|---|---|---|---|---|---|---|
| Religion | Base-BBQ | 0.6875 | 0.6833 | 0.6917 | 0.3167 | 0.1957 |
| | SAIC-BBQ | 0.6833 | 0.6833 | 0.6833 | 0.3167 | 0.1778 |
| | Base-vBBQ | 0.65 | 0.6667 | 0.6333 | 0.3333 | 0.2414 |
| | SAIC-vBBQ | 0.6458 | 0.6583 | 0.6333 | 0.325 | 0.1765 |
| SES | Base-BBQ | 0.9633 | 0.9712 | 0.9554 | 0.0089 | 0.0 |
| | SAIC-BBQ | 0.9648 | 0.9742 | 0.9554 | 0.0099 | 0.001 |
| | Base-vBBQ | 0.9504 | 0.9405 | 0.9603 | 0.0417 | 0.0062 |
| | SAIC-vBBQ | 0.9479 | 0.9435 | 0.9524 | 0.0347 | 0.0072 |
| Sexual | Base-BBQ | 0.6458 | 0.9896 | 0.3021 | -0.0104 | 0.0 |
| | SAIC-BBQ | 0.6458 | 0.9896 | 0.3021 | -0.0104 | 0.2258 |
| | Base-vBBQ | 0.7604 | 0.9375 | 0.5833 | 0.0208 | 0.0357 |
| | SAIC-vBBQ | 0.7448 | 0.9271 | 0.5625 | 0.0104 | -0.0182 |

# Yijun Wang

**Tonbridge School** （UK）
**Date of Birth**：**2009/04/28**

*With a strong foundation in coding, mathematics, and logical reasoning, I have excelled in competitive programming and algorithmic problem-solving, inspired by contests like the IOI. I particularly enjoy applying these skills across disciplines, connecting computer science with natural sciences and even classical studies. While studying in Britain, I discovered a deep love for Greco-Roman literature through the imaginative works of Vergil and Herodotus. The debate over Herodotus' reliability as a historian, together with my scientific curiosity, led me to explore the enduring problem of bias with modern innovations such as LLMs. Beyond academics, I am an avid chessplayer, history enthusiast, and enjoy sailing and playing tennis.*

## PROJECT EXPERIENCE

*TransScan - Scanning Latin dactylic verse using Transformers (IEEE Xplore & EI Conference Paper)*

*A Comparative Analysis of Shallow and Deep Learning Algorithms for Intrusion Detection Systems using NSL-KDD*

## AWARDS

*JACT Latin Camp 2025 (2 weeks of reading works by Cicero and Vergil)*

*British Informatics Olympiad 2025 (invited toIOIselection camp)*

*British Maths Olympiad (2024 Distinction + Bronze Medal)*

*British Maths Olympiad 2025, Round 2 qualification (top 100 in UK)*

*Cambridge Chemistry Challenge L6th (2024 Gold, 2025 Roentgenium)*

*UK Chemistry Olympiad, (2025 Gold), British Biology Olympiad, (2025 Gold)*

*PACT-Asia (Program in Algorithmic and Computational Thinking) Summer Program 2024, Best Student Award*

*SAT (Score: 1570/1600), TOEFL (Score :116/120)*