AesBiasBench: Evaluating Bias and Alignment in Multimodal Language Models for Personalized Image Aesthetic Assessment

Anonymous ACL submission

Abstract

Multimodal Large Language Models (MLLMs) are increasingly used in Personalized Image Aesthetic Assessment (PIAA), offering a scalable alternative to expert evaluation. However, their outputs may reflect subtle biases shaped by demographic cues such as gender, age, or education. In this work, we introduce AesBiasBench, a benchmark designed to evaluate MLLMs along two complementary axes: (1) the presence of stereotype bias, measured by how aesthetic evaluations vary across demographic groups; and (2) the alignment between model outputs and real human aesthetic preferences. Our benchmark spans three subtasks, Aesthetic Perception, Assessment, and Empathy, and introduces structured metrics (IFD, NRD, AAS) to quantify both bias and alignment. We evaluate 19 MLLMs, including proprietary models (e.g., GPT-4o, Claude-3.5-Sonnet) and open-source models (e.g., InternVL-2.5, Qwen2.5-VL). Results show that smaller models exhibit stronger stereotype bias, while larger models better align with human preferences. Adding identity information often amplifies bias, particularly in emotional judgment. These findings highlight the need for identity-aware evaluation frameworks for subjective vision-language tasks.

1 Introduction

004

005

011

012

017

035

040

043

Multimodal Large Language Models (MLLMs) have demonstrated impressive capabilities in vision-language tasks such as image recognition (Alayrac et al., 2022; Zhu et al., 2023), visual reasoning (Achiam et al., 2023; Wei et al., 2022), and visual question answering (Wu et al., 2023). Recently, these models have also been applied to Personalized Image Aesthetic Assessment (PIAA), which estimates the photographic or artistic quality of images based on individual preferences (Yang et al., 2022). PIAA applications include image retrieval, photo ranking, and creative recommendation (Ren et al., 2017).



Figure 1: Examples illustrate bias in the image aesthetic empathy task. (a) and (b) show stereotypical bias in model outputs that arise from inherited cognitive priors. (c) presents human preferences for the image, which serve as a reference for evaluating the alignment of model predictions with human judgments.

045

047

048

051

054

057

060

061

062

063

064

065

067

068

Despite their promise, MLLMs may exhibit aesthetic bias, systematic differences in output driven by demographic attributes such as gender, age, or education. Prior work has shown that even subtle biases in subjective tasks can lead to skewed outcomes (Zangwill, 2003; Dhamala et al., 2021; Tamkin et al., 2023). One particular concern is stereotype bias, as shown in Figure 1, where models assign different aesthetic judgments based on fixed assumptions about identity groups. Despite ongoing efforts to audit and debias deployed models for greater fairness (Guo et al., 2022; Smith et al., 2023; Dige et al., 2024; Li et al., 2024a,b), implicit and often-overlooked aesthetic biases continue to persist. Moreover, bias detection alone does not explain whether these deviations are problematic. Some output variation may simply reflect valid preference alignment with real human judgments. To address this, we complement bias measurement with an explicit evaluation of *alignment*, how closely model outputs match the aesthetic preferences of human users from corresponding demographic groups.

To support this dual analysis, we introduce **Aes-BiasBench**, a benchmark for assessing both stereo-

type bias and preference alignment in MLLMs applied to PIAA. Our benchmark covers three subtasks. The first, Aesthetic Perception, concerns the evaluation of low-level technical properties such as sharpness, lighting, and color. The second, Aesthetic Assessment, captures subjective evaluations of overall visual appeal and composition. The third, Aesthetic Empathy, targets the emotional impact conveyed or evoked by an image. For each subtask, we define dedicated metrics to quantify both bias and alignment, including Identity Frequency Disparity (IFD), Normalized Representation Disparity (NRD) and Aesthetic Alignment Score (AAS).

069

087

100

101

102

103

104

105

106

107

108

109

110

111

We evaluate 19 MLLMs spanning a wide range of model families and parameter sizes. The results show that smaller models tend to exhibit stronger stereotype bias, while larger models demonstrate both improved fairness and closer alignment with human preferences. In perception and assessment tasks, model outputs often align most closely with the preferences of female users aged 22 to 25 with a university education. In the empathy task, model responses align with female preferences by default, but shift toward male preferences when gender information is made explicit. This shift highlights strong sensitivity to identity cues rather than neutrality. By analyzing both bias and alignment, Aes-BiasBench enables a more complete understanding of fairness and demographic sensitivity in MLLMs. It provides a foundation for future work on socially aware and user-aligned multimodal systems.

The contributions of this work are threefold:

- Revealing stereotype biases in MLLMs for PIAA using tailored metrics that quantify group-specific deviations.
- Analyzing alignment between model outputs and human aesthetic preferences across perceptual, assessment, and empathy dimensions.
- Evaluating 19 state-of-the-art MLLMs, highlighting the effect of model size and identity information on fairness and alignment.

2 Related Work

2.1 Personalized Image Aesthetic Assessment

112Image aesthetic assessment (IAA) aims at compu-
tationally evaluating image quality based on pho-
tographic rules(Deng et al., 2017). Due to signifi-
cant variations in aesthetic preferences among in-
dividuals, image aesthetics can be categorized into

Generic Image Aesthetic Assessment (GIAA) and Personalized Image Aesthetic Assessment (PIAA). Regarding GIAA, early studies focused on designing and extracting image features, mapping them to annotated aesthetic labels. As a result, numerous IAA datasets have emerged to support research in this field (Dhar et al., 2011; Murray et al., 2012; Yi et al., 2023). 117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

166

167

Personalized Image Aesthetic Assessment aims to capture the unique aesthetic preferences of individuals (Yang et al., 2022). Existing methods typically rely on generic image aesthetic assessment (GIAA), incorporating rich attributes to facilitate specific aesthetic predictions (Li et al., 2020). Ren et al. (2017) found that individual aesthetic preferences have a strong correlation with image content and aesthetic attribute, and proposed residual scores to modify generic aesthetic scores into personalized aesthetic scores. Zhu et al. (2020) trained PIAA models by fine-tuning a pretrained GIAA model on personal rating data, while Cui et al. (2020) utilized GIAA model as feature extractor to capture deep feature representing aesthetic preferences. Instead based on priorledge of GIAA, Hou et al. (2022) directly estimated personalized aesthetic experiences by analyzing interaction matrices, which represents interaction between image content and user preferences. At the same time, the Q-instruct (Wu et al., 2024a) framework has improved MLLMs' low-level visual capabilities across multiple base models, and Q-align (Wu et al., 2023) leveraged the visual power of MLLMs to propose a new scoring method. These advances have laid the groundwork for using MLLMs in PIAA tasks.

2.2 Biases in MLLMs

The recent success of large language models (LLMs) has fueled exploration into visionlanguage interaction, leading to the emergence of multimodal large language models (MLLMs). These models have demonstrated strong capabilities in dialogue based on visual inputs. Given their advanced visual understanding, MLLMs can be leveraged to tackle various multimodal tasks related to high-level vision, including image aesthetic assessment(Zhou et al., 2024). However, the inherent biases in MLLMs may introduce systematic distortions in image evaluations, leading to biased aesthetic assessments.

Recent studies have explored the response biases in LLMs, which often influenced by various con-



Figure 2: AesBiasBench framework for stereotype bias measurement and aesthetic alignment evaluation. The model's default prompt includes an image < img > and task t, while the personalized prompt adds a demographic group g. After obtaining model responses for all images, IFD and NRD detect stereotype bias, while AAS identifies alignment, revealing the demographic group the model's aesthetic preferences align with.

168 textual and cultural factors(Gallegos et al., 2024; Tjuatja et al., 2023). Such biases also appear in 169 MLLMs, where visual and textual modalities can 170 interact in ways that reinforce existing societal bi-171 ases(Chen et al., 2024a). These biases are com-172 monly detected and analyzed through its manifes-173 tations in models outputs(Lin et al., 2024; Kumar 174 et al., 2024; Naous et al., 2023). Specifically, Jiang et al. (2024) revealed differences in occupations, 176 descriptions, and personality traits due to social 177 gender and racial biases across both visual and lan-178 guage modalities. Building on this line of work, 179 we focus on aesthetic biases that emerge when MLLMs evaluate images conditioned on identity in-181 formation. We first examine stereotype bias, where 182 models produce systematically different aesthetic 183 judgments across demographic groups. We then evaluate whether these biased outputs align with 185 the aesthetic preferences of corresponding human groups, highlighting how model behavior may reflect or distort real-world preferences.

3 Methodology

189

190

192

193

194

196

197

198

199

200

3.1 Preliminaries

This section introduces our definition and design of bias quantification when MLLMs applied to personalized image aesthetic assessment labeling. The bias quantification problem includes four components: the image needed to be assessed < img >, the specific assess task t, the identity group G and the MLLM used for quality assessment $M(\cdot)$. We can collect the response from the MLLM as follows:

$$M_t(<\operatorname{img}>|G=g) \tag{1}$$

where $M(\cdot) \in \Delta$ and Δ denotes the output format. Following (Huang et al., 2024), we focus on three assessment tasks t, including Aesthetic Perception which representing the perceived technical quality of the image, Aesthetic Assessment which representing the subjective aesthetic appeal of the image, and Aesthetic Empathy which capturing the emotional response evoked by the image. 201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

226

227

228

229

230

231

232

233

234

For identity group, we evaluate the bias across three different group categories, including age, gender and education group. We divide the individuals to different identities in each group category. For age group, we divide individuals into five different identities: 18 to 21, 22 to 25, 26 to 29, 30 to 34, and 35 to 40. For education group, we classify them into five education levels: junior high school, technical secondary school, senior high school, university, and junior college. For gender category, we consider male and female.

We define the output format Δ for each of the three tasks. For Aesthetic Perception and Aesthetic Assessment, $\Delta = \{\text{positive, normal, negative}\}$. For Aesthetic Empathy, $\Delta = \{\text{amusement, excitement, contentment, awe, disgust, sadness, fear, neutral}\}.$

3.2 Quantifying Bias

To analyze stereotype bias, we propose two metrics: Identity Frequency Disparity (IFD) and Normalized Representation Disparity (NRD). Identity Frequency Disparity (IFD) measures differences in how often the model assigns specific aesthetic evaluations Δ to various identity groups. This metric quantifies disparities in frequency, revealing potential biases in how different identities are assessed.

275 276

277

278

- 281 282
- 284
- 285 286
- 287
- 289
- 290
- 291

- 295

297

298 299

301

302

303 304

- 306 307 308
- 309
- 310
- 311
- 312

Normalized Representation Disparity (NRD) examines the model's preferences and emotional responses toward different types of images across identities. By normalizing for baseline differences in representation, NRD captures variations in the model's perceptions and affective reactions that may indicate bias. Together, these metrics provide a structured approach to identifying and quantifying stereotype bias in the model's behavior. They are defined as:

236

241

244

245

246

247

248

253

254

257

261

262

263

265

266

270

271

272

273

274

IFD(t) =
$$\frac{1}{|\Phi| \times n_{\Delta}} \sum_{k=1}^{n_{\Delta}} \sum_{G \in \Phi} \sum_{g=1}^{n_{G}} |p_{g,k} - p_{G,k}|,$$
(2)

where $p_{g,k}$ represents the proportion of choice k made by identity group g, n_G is the number of identity groups in category G, Φ is the set of all group categories and n_{Δ} is the number of the output choice.

The Normalized Representation Disparity (NRD) measures the disparities in the model's output $M(\cdot)$ between different identity groups g for a given task t, normalized by the total sentiment for each output $M(\cdot)$ across all identity groups. It is defined as:

$$\operatorname{NRD}(t) = \frac{1}{n_{\Delta}} \sum_{k=1}^{n_{\Delta}} \sqrt{\frac{1}{n_G} \sum_{g=1}^{n_G} \sum_{m \in \Omega} \left(F_{g,k}^m - \frac{1}{n_G} \right)^2},$$
(3)

$$F_{g,k}^{m} = \frac{S_{g,k}^{m}}{\sum_{h=1}^{n_{G}} S_{h,k}^{m}},$$
(4)

where $S_{g,k}^m$ is the number of times the outputs $M(\cdot) = k$ for identity group g and task t within image type m and Ω is the set of all image types.

3.3 Alignment Evaluation

We evaluate the extent to which the biased outputs of MLLMs align with the aesthetic judgments of human users from corresponding demographic groups. This analysis focuses on measuring how closely model outputs reflect real human preferences, providing a complementary perspective on the effects of stereotype bias. We conduct this evaluation from two perspectives:

• We examine which demographic groups the model's aesthetic judgments are more aligned with its default or pre-trained aesthetic preferences. This focuses on identifying whether the

model shows a stronger bias towards certain groups when no specific identity is specified.

• We explore which demographic groups the model's aesthetic judgments align more closely with human aesthetic preferences, when given the identity information explicitly. This helps identify whether the model's outputs reflect the actual preferences of different identity groups.

To measure the similarity between two outputs, we compute the similarity score using the Jensen-Shannon Divergence. Let O_q and O_h represent the model's outputs for images from groups g and hwhere $O_q, O_h \in \Delta$. To compute the JS divergence, we first map the discrete aesthetic choices in Δ to probability distributions using a one-hot encoding scheme, obtaining E_q and E_h . The JS divergence between E_q and E_h can then be calculated as:

$$JS(E_g || E_h) = \frac{1}{2} \left[D_{KL}(E_g || M) + D_{KL}(E_h || M) \right],$$
(5)

where M is the average distribution of E_g and $E_h, M = \frac{E_g + E_h}{2}$. And D_{KL} is the Kullback-Leibler (KL) Divergence, given by:

$$D_{\mathrm{KL}}(E \parallel M) = \sum_{j} E(j) \log\left(\frac{E(j)}{M(j)}\right).$$
(6)

To evaluate the alignment, we define the similarity score as:

$$S(g) = 1 - \mathbf{JS}(E_g \parallel E_h), \tag{7}$$

and the Aesthetic Alignment Score (AAS) is defined as follows:

$$AAS_t(g) = S_t(g) - \bar{S}_t, \qquad (8)$$

where $S_t(g)$ is the similarity score of the current identity in task t and \bar{S}_t is the mean similarity score of the category G in task t.

This metric is designed to compare the relative accuracy across different demographic groups, highlighting potential disparities in the model's ability to align with human aesthetic evaluations.

Experiments and Results 4

4.1 Experimental Setup

1

Dataset. In our experiments, we investigate bias 313 in three identity dimensions: gender, age, and ed-314 ucation. Each dimension is specifically chosen to 315



Figure 3: Left: IFD scores heatmap across diverse set of models. Right: Radar chart of IFD scores for InternVL-2.5 series models, showing variations by model size. A higher IFD indicates a greater degree of stereotype bias.

investigate societal biases in aesthetic perceptions toward the respective groups. We perform extensive testing on a well-established dataset for personalized image aesthetic assessment (Yang et al., 2022), the Personalized Image Aesthetics Database with Rich Attributes (PARA). PARA comprises 31,220 images annotated by 438 human raters with rich feature annotations. Built upon it, we generate three types of task evaluations for the 31,220 images: aesthetic perception, aesthetic assessment, and emotional perception. The three tasks are evaluated by IFD, NRD, AAS, and similarity score to examine both stereotype bias and aesthetic alignment bias. To align with human raters, who typically judge

316

317

319

321

323

328

329

331

335

339

340

341

344

347

based on discrete text-defined levels, we convert continuous scores in the PARA dataset into discrete rating levels. In AesBiasBench (ABB), this ensures consistency with the output formats Δ and facilitates fair comparisons between model outputs and human judgments. We adopt equidistant intervals to convert scores into rating levels as Wu et al. (2023), which is to uniformly divide the range between the highest score (M) and lowest score (m) into three distinct intervals.

$$L(s) = l_i \tag{9}$$

where $m + \frac{i-1}{3} \times (M-m) < s \le m + \frac{i}{3} \times (M-m)$.

Models. In this work, we investigate a diverse set of models, including InternVL2.5 (1B, 2B, 4B, 8B, 26B, 38B) (Chen et al., 2024b,c,d), Qwen2.5-VL (3B, 7B) (Yang et al., 2024), LLaVA (v1.5-7B, v1.6-vicuna-7B) (Liu et al., 2023a,b), LLaMA-**3.2** (11B-vision-instruct) (Grattafiori et al., 2024), mPLUG-Owl3 (7B) (Ye et al., 2024), Mono-InternVL (2B) (Luo et al., 2024), Phi-3.6-Vision (Abdin et al., 2024), GLM-4V (9B) (GLM et al., 2024), and DeepSeek (VL2) (Wu et al., 2024b). We also include closed-source models such as Claude-3.5-Sonnet, Gemini-2.0-flash and GPT-40 in our analysis. This comprehensive selection enables a systematic evaluation of biases across a wide range of architectures and scales. With this setup, we can compare bias variations within the same model series across different sizes. as well as across models of similar sizes. Such comparisons provide deeper insights into how model architecture, scale, and training paradigms influence bias.

348

349

350

351

354

355

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

379

4.2 Stereotype Bias Analysis

4.2.1 Existence of Bias in MLLMs

We quantify stereotype bias in MLLMs performing PIAA using two metrics: Identity Fairness Deviation (IFD) and Normalized Response Deviation (NRD). The heatmap in Figure 3 shows the IFD scores across multiple models, indicating substantial identity-related biases, where higher IFD values reflect stronger bias. Among these, the InternVL2.5 model series consistently shows lower IFD values, suggesting better fairness across demographic identities.

Additionally, Figure 4 (left) illustrates NRD scores, confirming strong biases, particularly evident in empathy-driven aesthetic tasks. Gender is



Figure 4: Left: NRD scores for age, gender, and education across three tasks. Right: $F_{g,k}^m$ scores of fear emotion from different groups for aesthetic empathy task in Claude-3.5-Sonnet, illustrating stereotype bias.

consistently identified as a major influencing factor, with notably higher NRD scores across all evaluated models. This emphasizes significant differences in the emotional perception of images among different demographic groups.

To further illustrate this, Figure 4 right provides a detailed example using Claude-3.5-Sonnet in the empathy task. The model predicts that younger individuals, those with lower education levels, and females are more likely to exhibit fear responses. These results suggest that advanced models encode systematic differences in emotional aesthetic judgment across demographic groups in emotional aesthetic judgment, reinforcing the presence of subtle but persistent stereotypical biases in MLLMs.

4.2.2 Impact of Model Size on Bias

The radar chart in Figure 3 right shows the IFD scores across the InternVL2.5 series. The results reveal a clear inverse relationship between model size and stereotype bias: as the model size increases from 1B to 38B, the IFD scores consistently decrease. InternVL2.5-1B shows the highest level of bias, followed by 2B, 4B, and 8B, with each larger model displaying progressively lower bias. The largest models, 26B and 38B, yield the most stable and fair outputs. This trend indicates that identityrelated bias decreases consistently with increasing model size. 401

402

403

404

405

406

407

This pattern is not limited to the InternVL2.5 408 series. Similar trends are observed in other model 409 families, where smaller variants consistently ex-410 hibit higher IFD scores than their larger counter-411 parts, indicating stronger stereotype bias. While 412 this may appear to reflect the effect of model ca-413 pacity alone, it is likely influenced by differences 414 in training data scale and diversity as well. Larger 415 models are often trained on broader and more bal-416 anced datasets, which may provide better coverage 417 of identity-related variations and contribute to more 418 equitable outputs. 419

400



Figure 5: AAS of the model on three tasks without identity information, showing the two most common identity patterns for each task. \circ , \diamond , and \triangle represent groups by gender, age, and education respectively.

4.3 Aesthetic Alignment Analysis

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

4.3.1 Default Aesthetic Preferences of Models

We begin by analyzing the default aesthetic alignment of MLLMs when no identity information is provided in the prompt. Using the Aesthetic Alignment Score (AAS), we measure the similarity between model outputs and the aesthetic preferences of different demographic groups across the three tasks.

The heatmap and radar plots in Figure 5 and summary statistics in Table 1 reveal clear and consistent demographic biases across tasks. All three tasks show a strong alignment with **female** aesthetic preferences, with 17 out of 19 models exhibiting this pattern. In terms of age, the **22–25** group dominates in Perception and Assessment, while Empathy shows a shift toward the younger 18–21 group. Educational alignment is more taskspecific. The most consistent pattern appears in the Assessment task, where nearly all models align with the same group: **female**, aged **22–25**, with a **university** education.

Task-specific patterns also emerge. As shown in Figure 5, the points in the radar plot for the Empathy task are more tightly clustered, indicating that the AAS values are generally lower compared to the other tasks. This aligns with the observation

	Perception	Assessment	Empathy
Gender	female (17)	female (17)	female (17)
Age	22_25 (12)	22_25 (17)	18_21 (8)
Education	Tech (7)	University (17)	Junior (7)

Table 1: The number of models exhibiting the highest AAS with different demographic groups across three tasks. The table summarizes results from 19 models.

in Figure 3, where the Empathy task also exhibits lower IFD values. Together, these results show that the default models are more fair in the Empathy task and exhibit weaker alignment with human aesthetic preferences. 447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

4.3.2 Sensitivity to Identity in Aesthetic Preferences

To further examine how identity information influences aesthetic alignment, we analyze the consistency of identity patterns across tasks after explicitly including demographic attributes in the prompts.

As shown in Figure 6 and summary statistics in Table 2, adding explicit identity information reduces the number of models that share the same dominant aesthetic pattern. This shift reflects that model outputs are sensitive to demographic descriptors, indicating the absence of neutral or identityinvariant behavior. It indicates that aesthetic out-

7



Figure 6: AAS of the model on three tasks with identity information, showing the two most common identity patterns for each task. \circ , \diamond , and \triangle represent groups by gender, age, and education respectively.

	Perception	Assessment	Empathy
Gender	female (15)	female (14)	male (17)
Age	22_25 (10)	22_25 (15)	30_34 (8)
Education	Junior (7)	University (10)	University (6)

Table 2: The number of models exhibiting the highest AAS with different demographic groups across three tasks when explicit identity attributes are provided. The table summarizes results from 19 models.

puts are systematically influenced by identity descriptors, revealing latent social biases in the models.

In particular, Table 2 shows a striking shift in the Empathy task: 17 models align with male identities, which is a complete reversal from the identityagnostic setting, where 17 models had aligned with female. Table 3 illustrates this bias sensitivity, showing increased alignment with male preferences when gender is added.

As shown in Table 3, most models show a greater increase in similarity to male preferences after gender is specified, indicating higher sensitivity to male identity. Instead of exposing more balanced behavior, the inclusion of gender information reveals stronger model bias, with responses becoming more aligned to male-associated aesthetic patterns-a deviation possibly reflecting differences in training data composition or architectural design.

5 Conclusion

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

490

491

492

493

This paper introduced AesBiasBench, a benchmark for evaluating biases in MLLMs on PIAA tasks. To quantify stereotype bias, we proposed 488 two metrics: IFD and NRD. In addition, we used 489 the AAS to measure how model outputs correspond to human aesthetic preferences across demographic groups. Key findings include: (1) Stereotype bias is prevalent across models, with smaller models showing more pronounced deviations and larger 494

Model	$\Delta S_E(\mathbf{M})$	$\Delta S_E(\mathbf{F})$	Δ
DeepSeek-VL2	-0.0535	-0.0749	0.0214
Llama-3.2-11B-Vision	0.0074	-0.0021	0.0095
GPT-40	0.0395	-0.0748	0.1143
Phi-3.5-Vision	-0.0113	-0.0293	0.0180
Claude-3.5-Sonnet	0.1180	-0.1166	0.2346
Gemini-2.0-Flash	0.0274	-0.3780	0.4054
GLM-4V-9B	-0.0743	-0.1015	0.0272
mPLUG_Owl3	-0.0047	-0.0226	0.0179
Qwen2.5-VL-3B	0.0330	0.0196	0.0134
Qwen2.5-VL-7B	0.0287	0.0198	0.0089
InternVL2.5-1B	0.0220	-0.0244	0.0464
InternVL2.5-2B	0.0187	-0.1971	0.2158
InternVL2.5-4B	0.0085	-0.0022	0.0107
InternVL2.5-8B	-0.0007	-0.0126	0.0119
InternVL2.5-26B	-0.0160	-0.0324	0.0164

Table 3: $\Delta S_E(M)$ and $\Delta S_E(F)$ denote the changes in similarity scores to male and female aesthetic preferences, respectively, after adding gender identity to the prompt in the empathy task. Δ represents the incremental gain of male over female, computed as $\Delta S_E(M)$ $-\Delta S_E(F)$. The top 3 highest and lowest Δ values are highlighted using soft red and blue gradients.

models exhibiting lower IFD and NRD scores, indicating increased fairness with scale. (2) Model outputs align disproportionately with certain demographic groups, notably, female individuals aged 22-25 with university education-even when identity information is not provided. (3) Adding identity descriptors amplifies existing biases, as shown in the empathy task where alignment shifts more strongly toward male preferences, revealing heightened sensitivity to demographic cues rather than neutrality. These results highlight the importance of identity-aware evaluation and point to the need for fairness-oriented design in future MLLMs used for subjective and socially-influenced tasks.

limitations

509

515

516

517

519 520

522

523

524

525

526

529

530

532

533

534

535

536

537

540 541

542

545

546

547

548

550

552

554

555

558

510 Our AesBiasBench evluates existing MLLMs 511 along two complementary axes: (1) stereotype bias 512 and (2) human preference alignment. To make the 513 results more reliable, we indentify two possible 514 limitations:

> First, the analysis is restricted to three identity attributes: age, gender, and education. While these dimensions capture important aspects of demographic variation, other factors, such as culture, race, and religion, may also influence aesthetic preferences and model behavior. Incorporating a broader range of identity dimensions could enable a more comprehensive understanding of demographic bias in MLLMs.

Second, we evaluate 19 MLLMs, including proprietary models (e.g., GPT-40, Claude-3.5-Sonnet) and open-source models (e.g., InternVL2.5 and Qwen2.5-VL series). While this selection spans a range of model families and sizes, future work could explore a broader set of architectures, training strategies, and deployment contexts, which may reveal additional forms of bias or alternative alignment.

Ethics

In this study, we constructed AesBiasBench using the publicly accessible Personalized Image Aesthetics Database with Rich Attributes (PARA). No original data collection was conducted; all analyses relied solely on pre-existing dataset resources. To the best of our knowledge, the PARA dataset was developed in strict adherence to academic and scientific data collection protocols, ensuring compliance with ethical standards for research involving human subjects.

Our research does not involve any personally identifiable information (PII) or process private/sensitive user data. The demographic attributes utilized (e.g., age groups, gender, education levels) are provided in the PARA dataset as anonymized and aggregated metadata, with no individual-level data accessible. This design ensures that no participant can be re-identified through the study's analyses.

The core objective of this research is to systematically uncover and characterize biased behaviors of multimodal large language models (MLLMs) in personalized aesthetic judgment tasks. By quantifying demographic disparities in model outputs, we aim to foster greater awareness within the research community and contribute to the development of559more equitable, transparent, and socially account-
able AI systems. Our work aligns with the broader560ethical imperative to promote fairness in machine562learning, particularly in applications impacting hu-
man values and societal norms.564

References

565

566

567

568

569

575

576

584

591

592

593

595

597

598

599

610

611

612

613

614

615

616

617

618

619

621

- Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, and 1 others. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, and 1 others. 2022. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716– 23736.
- Dongping Chen, Ruoxi Chen, Shilin Zhang, Yinuo Liu, Yaochen Wang, Huichi Zhou, Qihui Zhang, Yao Wan, Pan Zhou, and Lichao Sun. 2024a. Mllmas-a-judge: Assessing multimodal llm-as-a-judge with vision-language benchmark. *arXiv preprint arXiv:2402.04788*.
- Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, and 1 others. 2024b. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*.
- Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, and 1 others. 2024c. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*.
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, and 1 others. 2024d. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24185–24198.
- Chaoran Cui, Wenya Yang, Cheng Shi, Meng Wang, Xiushan Nie, and Yilong Yin. 2020. Personalized image quality assessment with social-sensed aesthetic preference. *Information Sciences*, 512:780–794.
- Yubin Deng, Chen Change Loy, and Xiaoou Tang. 2017. Image aesthetic assessment: An experimental survey. *IEEE Signal Processing Magazine*, 34(4):80–106.
- J Dhamala, T Sun, V Kumar, S Krishna, Y Pruksachatkun, K-W Chang, and R Gupta. 2021. Bold: Dataset and metrics for measuring biases in openended language generation. In *Proceedings of the* 2021 ACM conference on Fairness, Accountability, and Transparency, pages 862–872.

Sagnik Dhar, Vicente Ordonez, and Tamara L Berg. 2011. High level describable attributes for predicting aesthetics and interestingness. In *CVPR 2011*, pages 1657–1664. IEEE.

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

- Omkar Dige, Diljot Singh, Tsz Fung Yau, Qixuan Zhang, Borna Bolandraftar, Xiaodan Zhu, and Faiza Khan Khattak. 2024. Mitigating social biases in language models through unlearning. *arXiv* preprint arXiv:2406.13551.
- Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. 2024. Bias and fairness in large language models: A survey. *Computational Linguistics*, pages 1–79.
- Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui, Jie Tang, Jing Zhang, Juanzi Li, and 37 others. 2024. Chatglm: A family of large language models from glm-130b to glm-4 all tools. *Preprint*, arXiv:2406.12793.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The Ilama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Yue Guo, Yi Yang, and Ahmed Abbasi. 2022. Autodebias: Debiasing masked language models with automated biased prompts. In *Proceedings of the* 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1012–1023.
- Jingwen Hou, Weisi Lin, Guanghui Yue, Weide Liu, and Baoquan Zhao. 2022. Interaction-matrix based personalized image aesthetics assessment. *IEEE Transactions on Multimedia*, 25:5263–5278.
- Yipo Huang, Quan Yuan, Xiangfei Sheng, Zhichao Yang, Haoning Wu, Pengfei Chen, Yuzhe Yang, Leida Li, and Weisi Lin. 2024. Aesbench: An expert benchmark for multimodal large language models on image aesthetics perception. *arXiv preprint arXiv:2401.08276*.
- Yukun Jiang, Zheng Li, Xinyue Shen, Yugeng Liu, Michael Backes, and Yang Zhang. 2024. Modscan: Measuring stereotypical bias in large vision-language models from vision and language modalities. *arXiv preprint arXiv:2410.06967*.
- Abhishek Kumar, Sarfaroz Yunusov, and Ali Emami. 2024. Subtle biases need subtler measures: Dual metrics for evaluating representative and affinity bias in large language models. *arXiv preprint arXiv:2405.14555*.
- Cheng Li, Mengzhuo Chen, Jindong Wang, Sunayana Sitaram, and Xing Xie. 2024a. Culturellm: Incorporating cultural differences into large language models. *arXiv preprint arXiv:2402.10946*.

778

780

781

782

783

784

785

786

732

Cheng Li, Damien Teney, Linyi Yang, Jindong Wang, Qingsong Wen, Xing Xie, and Jindong Wang. 2024b. Culturepark: Boosting cross - cultural understanding in large language models. *arXiv preprint arXiv:2405.15145*.

678

679

682

688

697

703

705

706

707

710

711

712

713

714

716

717

719

720

721

725

727

728

731

- Leida Li, Hancheng Zhu, Sicheng Zhao, Guiguang Ding, and Weisi Lin. 2020. Personality-assisted multi-task learning for generic and personalized image aesthetics assessment. *IEEE Transactions on Image Processing*, 29:3898–3910.
- Luyang Lin, Lingzhi Wang, Jinsong Guo, and Kam-Fai Wong. 2024. Investigating bias in llm-based bias detection: Disparities between llms and human perception. *arXiv preprint arXiv:2403.14896*.
 - Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023a. Improved baselines with visual instruction tuning.
 - Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. Visual instruction tuning. In *NeurIPS*.
- Gen Luo, Xue Yang, Wenhan Dou, Zhaokai Wang, Jiawen Liu, Jifeng Dai, Yu Qiao, and Xizhou Zhu. 2024. Mono-internvl: Pushing the boundaries of monolithic multimodal large language models with endogenous visual pre-training. *arXiv preprint arXiv:2410.08202*.
- Naila Murray, Luca Marchesotti, and Florent Perronnin. 2012. Ava: A large-scale database for aesthetic visual analysis. In 2012 IEEE conference on computer vision and pattern recognition, pages 2408–2415. IEEE.
- Tarek Naous, Michael J Ryan, Alan Ritter, and Wei Xu. 2023. Having beer after prayer? measuring cultural bias in large language models. arXiv preprint arXiv:2305.14456.
- Jian Ren, Xiaohui Shen, Zhe Lin, Radomir Mech, and David J Foran. 2017. Personalized image aesthetics. In Proceedings of the IEEE international conference on computer vision, pages 638–647.
- Brandon Smith, Miguel Farinha, Siobhan Mackenzie Hall, Hannah Rose Kirk, Aleksandar Shtedritski, and Max Bain. 2023. Balancing the picture: Debiasing vision-language datasets with synthetic contrast sets. *arXiv preprint arXiv:2305.15407*.
- A Tamkin, A Askell, L Lovitt, E Durmus, N Joseph, S Kravec, K Nguyen, J Kaplan, and D Ganguli. 2023.
 Evaluating and mitigating discrimination in language model decisions. *arXiv preprint arXiv:2312.03689*.
- Lindia Tjuatja, Valerie Chen, Sherry Tongshuang Wu, Ameet Talwalkar, and Graham Neubig. 2023. Do Ilms exhibit human-like response biases? a case study in survey design. *arXiv preprint arXiv:2311.04076*.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, and

1 others. 2022. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*.

- Haoning Wu, Zicheng Zhang, Erli Zhang, Chaofeng Chen, Liang Liao, Annan Wang, Kaixin Xu, Chunyi Li, Jingwen Hou, Guangtao Zhai, and 1 others. 2024a.
 Q-instruct: Improving low-level visual abilities for multi-modality foundation models. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 25490–25500.
- Haoning Wu, Zicheng Zhang, Weixia Zhang, Chaofeng Chen, Liang Liao, Chunyi Li, Yixuan Gao, Annan Wang, Erli Zhang, Wenxiu Sun, and 1 others. 2023. Q-align: Teaching lmms for visual scoring via discrete text-defined levels. arXiv preprint arXiv:2312.17090.
- Zhiyu Wu, Xiaokang Chen, Zizheng Pan, Xingchao Liu, Wen Liu, Damai Dai, Huazuo Gao, Yiyang Ma, Chengyue Wu, Bingxuan Wang, and 1 others. 2024b. Deepseek-vl2: Mixture-of-experts visionlanguage models for advanced multimodal understanding. arXiv preprint arXiv:2412.10302.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, and 22 others. 2024. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*.
- Yuzhe Yang, Liwu Xu, Leida Li, Nan Qie, Yaqian Li, Peng Zhang, and Yandong Guo. 2022. Personalized image aesthetics assessment with rich attributes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19861– 19869.
- Jiabo Ye, Haiyang Xu, Haowei Liu, Anwen Hu, Ming Yan, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. 2024. mplug-owl3: Towards long image-sequence understanding in multi-modal large language models. *Preprint*, arXiv:2408.04840.
- Ran Yi, Haoyuan Tian, Zhihao Gu, Yu-Kun Lai, and Paul L Rosin. 2023. Towards artistic image aesthetics assessment: a large-scale dataset and a new method. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22388– 22397.

Nick Zangwill. 2003. Aesthetic judgment.

- Zhaokun Zhou, Qiulin Wang, Bin Lin, Yiwei Su, Rui Chen, Xin Tao, Amin Zheng, Li Yuan, Pengfei Wan, and Di Zhang. 2024. Uniaa: A unified multi-modal image aesthetic assessment baseline and benchmark. *arXiv preprint arXiv:2404.09619*.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*.

Hancheng Zhu, Leida Li, Jinjian Wu, Sicheng Zhao,
Guiguang Ding, and Guangming Shi. 2020. Personalized image aesthetics assessment via meta-learning
with bilevel gradient optimization. *IEEE Transac- tions on Cybernetics*, 52(3):1798–1811.