

When Cultures Meet: Multicultural Text-to-Image Generation

Anonymous ACL submission

Abstract

Text-to-image generation models have achieved strong performance in culturally homogeneous settings, yet their ability to generate *multicultural scenes*—where people and landmarks originate from different cultures—remains largely unexplored. We introduce *multicultural text-to-image generation* as a new task and present the first benchmark designed to study this setting. Our dataset contains 9,000 images spanning five countries, three age groups, two genders, 25 historical landmarks, and five languages. Using this benchmark, we analyze the behavior of state-of-the-art text-to-image models across multiple dimensions, including alignment, image quality, aesthetics, knowledge, and fairness. As one strategy for composing cultural and demographic information, we explore MosAIG, a Multi-Agent framework that enhances multicultural Image Generation by leveraging LLMs with distinct cultural personas. Our analysis shows that richer prompt composition can improve image quality and cultural grounding compared to simple prompts, while revealing substantial disparities across languages and demographic groups. We release our dataset and code at <https://anonymous.4open.science/r/MosAIG>.

1 Introduction

Societies worldwide are increasingly diverse, shaped by global travel and migration (Castles et al., 2103). This multicultural reality poses important challenges for Artificial Intelligence (AI), where robust representation of diverse populations is essential for equity and inclusivity (Hershcovich et al., 2022; Naous et al., 2024; Mihalcea et al., 2025). However, most datasets used for text-to-image generation focus on narrow demographics—predominantly Western, adult, and male—and largely depict single-culture scenarios (e.g., *a Chinese temple, an Indian market*) (Liu et al., 2024; Kannan et al., 2024). Such representations fail to capture common multicultural

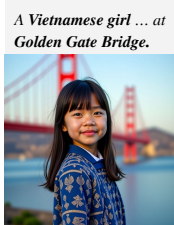
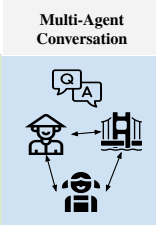

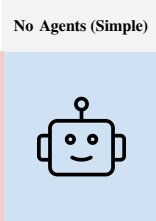
	Ours	Multi-Agent Conversation
	Inclusive Demographics (age, gender, country, language) Multicultural Images (landmark - person)	
	Related Work Biased Demographics (adult, male, Western, English) One Culture per Image (United States)	No Agents (Simple) 

Figure 1: Most existing datasets emphasize singular cultural contexts (e.g., the Golden Gate Bridge depicted primarily with American visitors or as a standalone monument). In contrast, real-world scenes often involve people from different cultural backgrounds sharing spaces and experiences. Modeling such multicultural interactions enables richer and more realistic image generation.

interactions, for example *a Chinese girl visiting the Golden Gate Bridge*, limiting the applicability of text-to-image models in real-world, culturally diverse settings (Hershcovich et al., 2022; Bhatia et al., 2024).

Despite recent efforts to evaluate cultural awareness in vision–language models, existing benchmarks and analyses primarily consider one culture per image. To date, there has been no systematic study of *multicultural text-to-image generation*, where elements from different cultural origins—such as a person and a landmark—co-exist within the same visual scene. We introduce this setting as a new task and study how current text-to-image models handle cultural composition, demographic variation, and cross-cultural grounding.

Specifically, we examine two key dimensions: (1) the demographic attributes of the depicted person, and (2) the multicultural interaction between the person and the landmark (e.g., the

Golden Gate Bridge). We consider four demographic aspects—age, gender, nationality, and language—together with cross-cultural landmarks (Figure 1). By systematically exploring these factors and their intersections, we aim to better understand the strengths and limitations of state-of-the-art text-to-image models in multicultural settings.

Our work is guided by the following research questions:

RQ1: How accurately do state-of-the-art text-to-image models depict people from one culture within the context of a landmark associated with a different culture?

RQ2: How does text-to-image generation performance vary across demographic groups and languages?

RQ3: What modeling strategies help improve multicultural text-to-image generation?

Contributions. **First, we release the first dataset of 9,000 images designed to study multicultural interactions**, depicting people and landmarks from different cultures across five countries, three age groups, two genders, 25 historical landmarks, and five languages. **Second, we explore MosAIG, a multi-agent prompting framework that decomposes cultural and demographic aspects** during caption construction as an effective strategy for addressing the task. **Finally, we analyze multicultural text-to-image generation through automated and human evaluation, revealing demographic and linguistic disparities.**

2 Related Work

Cultural Evaluation in Language and Vision Models. Recent work has made substantial progress in modeling and evaluating cultural awareness in language models through large multilingual benchmarks (Pawar et al., 2024; Romanou et al., 2025; Singh et al., 2025). In the vision–language domain, benchmarks such as CVQA (Romero et al., 2024) and GlobalRG (Bhatia et al., 2024) evaluate culturally grounded question answering, retrieval, and visual grounding. While multi-agent approaches have been explored for cross-cultural reasoning in multimodal systems (Guo et al., 2024; Han et al., 2024), prior work such as MosAIC (Bai et al., 2025) focuses on image captioning in single-culture settings, where models describe visual content post hoc. In contrast, we study *text-to-image generation in multicultural settings*, where models

must jointly reason about and synthesize multiple cultural, demographic, and landmark-specific cues into a single coherent visual scene. This setting constitutes a particularly challenging test of cultural competence, as failures in grounding, representation, or composition are directly reflected in the generated images.

Text-to-Image Generation Models and Benchmarks. Text-to-image generation has advanced rapidly with models such as Stable Diffusion-XL (Podell et al., 2023), DALL-E 3 (Betker et al., 2023), and FLUX (Labs, 2024). Agentic approaches like GenArtist (Wang et al., 2024) focus on unified generation and editing pipelines, whereas our work emphasizes multicultural and multilingual evaluation rather than model design. Existing benchmarks, including TIFA (Hu et al., 2023), GenEval (Ghosh et al., 2024), and GenAIBench (Lin et al., 2025), primarily assess technical properties such as realism, faithfulness, and compositionality. More recent efforts, such as HEIM (Lee et al., 2024), incorporate socially situated dimensions including bias, toxicity, and aesthetics (Hartwig et al., 2024), but do not explicitly address multicultural scene composition.

Cultural and Linguistic Gaps in Text-to-Image Generation. Despite these advances, most text-to-image systems and evaluations remain centered on a limited set of high-resource languages, leaving many linguistic communities underserved. While models such as Taiyi-Diffusion-XL (Wu et al., 2024) and AltDiffusion (Ye et al., 2024) expand multilingual input coverage, a broader cultural gap persists (Liu et al., 2024), as existing benchmarks rarely capture cross-cultural interactions or multicultural contexts (Hershcovich et al., 2022; Mihalcea et al., 2025; Saha et al., 2025).

Data Diversity and Cultural Competence. Recent work has begun to assess cultural competence in text-to-image generation. For example, CUBE (Kannen et al., 2024) and TIFA (Hu et al., 2023) evaluate cultural awareness and diversity, but remains limited to single-culture depictions per image. To our knowledge, no prior work systematically studies *multicultural image generation*, where multiple cultures co-exist within the same scene. Our work addresses this gap by introducing a benchmark and analysis framework for multicultural text-to-image generation.

3 Multicultural Image Generation

Culture is a multifaceted concept, meaning different things to different people at different times (Adilazuarda et al., 2024). In this work, we adopt the definition of Nguyen et al. (2023) and focus on visual cultural elements, such as clothing and historical landmarks.

We introduce *multicultural image generation* as a new task that evaluates how text-to-image models represent elements from multiple cultures within a single image—specifically, a person from one cultural background depicted in the context of a landmark from another. In addition to cultural origin, we examine demographic attributes and their intersections, including age, gender, and language¹. To address this task, we introduce MosAIG, a novel framework for Multi-Agent Image Generation, as illustrated in Figure 2. Our framework generates comprehensive image captions that are used to generate more accurate multicultural images using off-the-shelf image generation models. This framework is built around a multi-agent interaction model, as described below.

3.1 Multi-Agent Interaction Model

We introduce a multi-agent setup to emulate collaboration between demographically diverse groups. Our setup contains five agents, with specific roles: one Moderator Agent, three Social Agents, and one Summarizer Agent, as illustrated in Figure 2.

Moderator Agent. The Moderator Agent obtains demographic (age, gender, nationality) information about the person, the name of the landmark (e.g., Taj Mahal), and the language of the caption as input. The Moderator Agent then assigns tasks to the Social agents, instructing them to focus on the visually relevant aspects of the input information.

Social Agents. The Social Agents interact by asking each other relevant questions to create an image caption according to the information provided by the Moderator Agent. Each Social Agent assumes a *persona*: the first agent represents the culture of the person in the image, the second agent represents the age and gender of the person, and the last agent represents the historical landmark. Each agent generates an initial description of their persona. Then, by interacting through multiple rounds of question-answering conversations, each agent creates a more comprehensive image description.

¹All demographics are listed in Appendix Table 1.

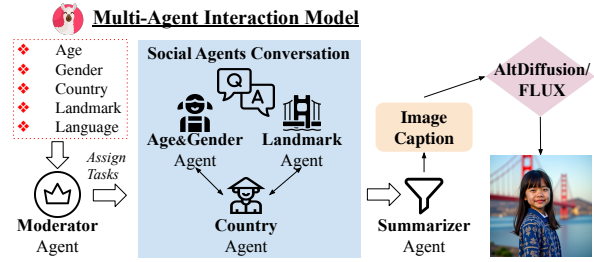


Figure 2: Overview of MosAIG, our framework for Multi-Agent Image Generation. The framework includes a multi-agent interaction model that generates an image caption from demographic information (person age, gender, country, landmark, and caption language), which is then used by an image generation model to create a multicultural image of a landmark and a person.

Summarizer Agent. The Summarizer Agent collects the three descriptions from the Social Agents and summarizes them into a final image caption with a maximum length of 77 tokens.

Social Agents Conversation. At the start, the three Social Agents—Country Agent, Landmark Agent, and Age-Gender Agent—receive demographic information and tasks from the Moderator Agent. The Country Agent processes nationality information and describes traditional attire, which is then evaluated by the Age-Gender Agent (e.g., “Is this attire suitable for a young female?”). Adjustments, such as modifying the color or style of a garment to suit the individual’s age, are made accordingly. The Landmark Agent describes the landmark architecture, and its descriptions are refined based on feedback from the Country Agent (e.g., “How do Vietnamese visitors typically interact with this landmark?”), ensuring cultural authenticity. The Age-Gender Agent generates demographic descriptions, which are cross-checked with the Country Agent to ensure culturally appropriate accessories and mannerisms. After two rounds of conversation, the agents enhance and refine the descriptions with culturally sensitive and contextually rich details. Once the iterative improvement process is complete, the refined descriptions are passed to the Summarizer Agent, which condenses them into a final 77-token prompt capturing the cultural and contextual nuances. The prompts and implementation details are provided in Figure 8 and Appendix C.

3.2 Image Generation Models

We evaluate our generated image captions using two different state-of-the-art image generation models: AltDiffusion (Ye et al., 2024) and

246 FLUX (Labs, 2024).

247 **AltDiffusion.** AltDiffusion² (Ye et al., 2024) is
248 one of the very few multilingual open-source image
249 generation models. The model aligns multilingual
250 language models with diffusion models to gener-
251 ate high-quality images from text across multiple
252 languages. The model builds on CLIP (Radford
253 et al., 2021), replacing its text encoder with XLM-
254 R (Conneau et al., 2020) and employing a two-stage
255 training process that combines teacher learning and
256 contrastive learning. AltDiffusion supports 18 dif-
257 ferent languages; we select five—English, German,
258 Hindi, Spanish, and Vietnamese—based on the an-
259 notators’ expertise. The model processes text inputs
260 with a maximum length of 77 tokens.

261 **FLUX.** FLUX.1-dev³ (Labs, 2024) is a state-of-
262 the-art, widely used, open-source text-to-image
263 model designed for English-language prompts.
264 Due to computational constraints, we employ
265 Flux.1 Lite⁴ (Daniel Verdú, 2024), an 8B-
266 parameter transformer model, more efficient vari-
267 ant distilled from FLUX.1-dev.

268 3.3 Simple vs. Multi-Agent Image Generation

269 Simple models generate images based on prede-
270 fined captions, whereas multi-agent models uti-
271 lize dynamically generated captions derived from
272 multi-agent interactions. For instance, when pro-
273 vided with demographic details such as “Viet-
274 namese” (nationality), “child” (age), “female”
275 (gender), “Golden Gate Bridge” (landmark), and
276 “English” (caption language), the resulting image
277 captions differ between the two approaches. Multi-
278 agent models generate captions that provide richer
279 contextual information, including detailed descrip-
280 tions of the landmark’s architecture and surround-
281 ings, as well as a more nuanced depiction of the
282 person’s appearance, particularly focusing on cloth-
283 ing and facial features, as shown below⁵.

284 **Simple caption:** *A Vietnamese girl wearing traditional attire,*
285 *standing in front of the Golden Gate Bridge.*

286 **Multi-agent caption:** *A 12-year-old Vietnamese girl in Áo*
287 *Dài, standing on the Golden Gate Bridge, with the San*
288 *Francisco Bay’s blue waters and the bridge’s orange-red*
289 *towers in the background.*

²<https://huggingface.co/BAAI/AltDiffusion-m18>

³<https://huggingface.co/black-forest-labs/FLUX.1-dev>

⁴<https://huggingface.co/Freepik/flux.1-lite-8B-alpha>

1-1ite-8B-alpha

⁵All the captions are shown in our code repository.

290 4 Evaluation and Results

291 We employ both automated metrics and human
292 evaluation to provide a holistic and comprehensive
293 assessment of the generated images.

294 4.1 Evaluation Metrics

295 We adopt a set of automated evaluation metrics that
296 assess text-to-image generation along five comple-
297 mentary dimensions: **Alignment**, **Quality**, **Aes-**
298 **thetics**, **Knowledge**, and **Fairness**. Together, these
299 metrics capture both technical properties of gener-
300 ation—such as semantic correspondence and vi-
301 sual fidelity—as well as socially situated aspects,
302 including representational consistency across de-
303 mographic groups (Lee et al., 2024). Given the
304 known limitations of any single automatic metric,
305 we combine multiple evaluators and complement
306 them with human judgment.

307 **Alignment.** We measure text–image alignment
308 using CLIPScore (Hessel et al., 2021), a widely
309 adopted, reference-free metric that computes co-
310 sine similarity between joint text and image em-
311 beddings and enables scalable evaluation. While
312 CLIPScore provides a useful proxy for semantic
313 correspondence, it does not capture all aspects of
314 visual grounding or compositional correctness. Ac-
315 cordingly, we interpret alignment scores cautiously
316 and complement them with human evaluation, and
317 we encourage future work to incorporate image-
318 based classifiers for more direct assessment of vi-
319 sual attribute realization (see Limitations). CLIP-
320 Score values range from -1 to $+1$, with higher
321 values indicating stronger alignment.

322 **Quality.** We assess image quality using the In-
323 ception Score (IS) (Salimans et al., 2016), which
324 evaluates both visual fidelity and output diversity
325 based on predictions from an Inception-v3 classi-
326 fier. Lower scores typically reflect poor realism
327 or limited variation, while higher scores indicate
328 more realistic and diverse images. Although IS
329 does not directly assess semantic correctness, it
330 provides a complementary signal for overall visual
331 plausibility.

332 **Aesthetics.** Aesthetic quality captures visual ap-
333 peal beyond semantic correctness, including sharp-
334 ness, color harmony, composition, and overall clar-
335 ity. We use a SigLIP-based aesthetic predictor⁶,
336 which assigns scores on a 1–10 scale. This met-
337 ric prioritizes perceptual attributes and may be less
338 sensitive to semantic or cultural accuracy, making it

⁶<https://github.com/discus0434/aesthetic-predictor-v2-5>

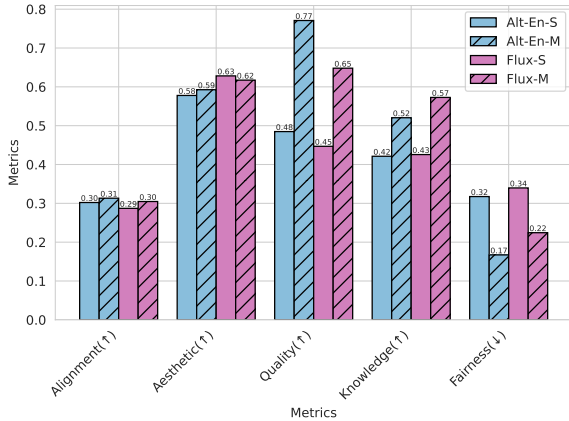


Figure 3: Multi-agent prompting improves *Quality*, *Knowledge*, and *Fairness* relative to simple prompts, while achieving comparable performance in *Alignment* and *Aesthetics*. Scores are normalized to [0–1]; higher is better except for Fairness.

particularly informative when interpreted alongside alignment and knowledge metrics.

Fairness. We evaluate fairness as the consistency of model performance across demographic substitutions. Following prior work, we modify captions by changing demographic attributes such as *gender*, *age*, or *nationality*, while keeping all other elements fixed. Given an original caption–image pair (c, I) , we construct a modified caption c' (e.g., replacing “boy” with “girl”) and generate a corresponding image I' . Fairness is measured as the absolute difference in alignment:

$$\Delta S = |S(c, I) - S(c', I')|.$$

Lower values of ΔS indicate more consistent behavior across demographic groups, while higher values suggest potential representational disparities. This metric captures relative changes rather than absolute bias and is best interpreted comparatively across models.

Knowledge. We assess world knowledge by evaluating a model’s sensitivity to landmark identity. Given a caption–image pair (c, I) , we replace the referenced historical landmark in the caption while keeping the image fixed, yielding (c', I) . We compute:

$$\Delta S = S(c, I) - S(c', I).$$

A model with stronger landmark knowledge should exhibit larger alignment differences when the caption references an incorrect landmark.

4.2 Multi-Agent Interaction Results

Figure 3 compares multi-agent and simple prompting strategies across five evaluation dimensions.

Overall, multi-agent prompting yields consistent improvements in **Quality**, **Knowledge**, and **Fairness**, while achieving comparable performance in **Alignment** and **Aesthetics**.

The largest gains are observed in **Quality**. Multi-agent models achieve substantially higher scores than simple prompts (0.77 vs. 0.48 for Alt-En and 0.65 vs. 0.45 for Flux-En), indicating more visually coherent and photorealistic outputs. We attribute these improvements to richer prompt composition that more explicitly specifies culturally grounded visual details, which appears to reduce under-specified or visually inconsistent generations. Across both prompting strategies, Alt consistently attains higher Quality scores than Flux, likely reflecting differences in background sharpness and image fidelity between the underlying generation models. In contrast, **Aesthetic** scores remain largely unchanged. This suggests that multi-agent prompting primarily affects semantic and compositional correctness rather than stylistic attributes emphasized by aesthetic predictors, which tend to prioritize surface-level visual appeal.

Multi-agent prompting also improves **Knowledge** and **Fairness**. Knowledge scores increase from 0.42 to 0.52 for Alt-En and from 0.43 to 0.57 for Flux-En, indicating stronger sensitivity to landmark-specific information when cultural context is more explicitly specified. Fairness scores—where lower values indicate smaller performance disparities—are substantially reduced (0.17 vs. 0.32 for Alt-En and 0.22 vs. 0.34 for Flux-En), suggesting more consistent behavior across demographic substitutions. These results indicate that decomposing cultural and demographic cues during prompt construction can mitigate uneven performance across social groups.

Improvements in **Alignment** are more modest and not statistically significant in aggregate. However, disaggregated analysis reveals consistent gains across several demographic dimensions, including *adults* (0.30 vs. 0.27), *females* (0.31 vs. 0.28), and multiple countries such as *Germany*, *India*, and *Vietnam* (see Appendix E.1). This pattern suggests that richer cultural specification can improve semantic correspondence for particular population groups, even when overall alignment scores remain similar.

4.3 Ablation Studies

We also perform ablation studies to assess MoSAIG’s performance across demographics.

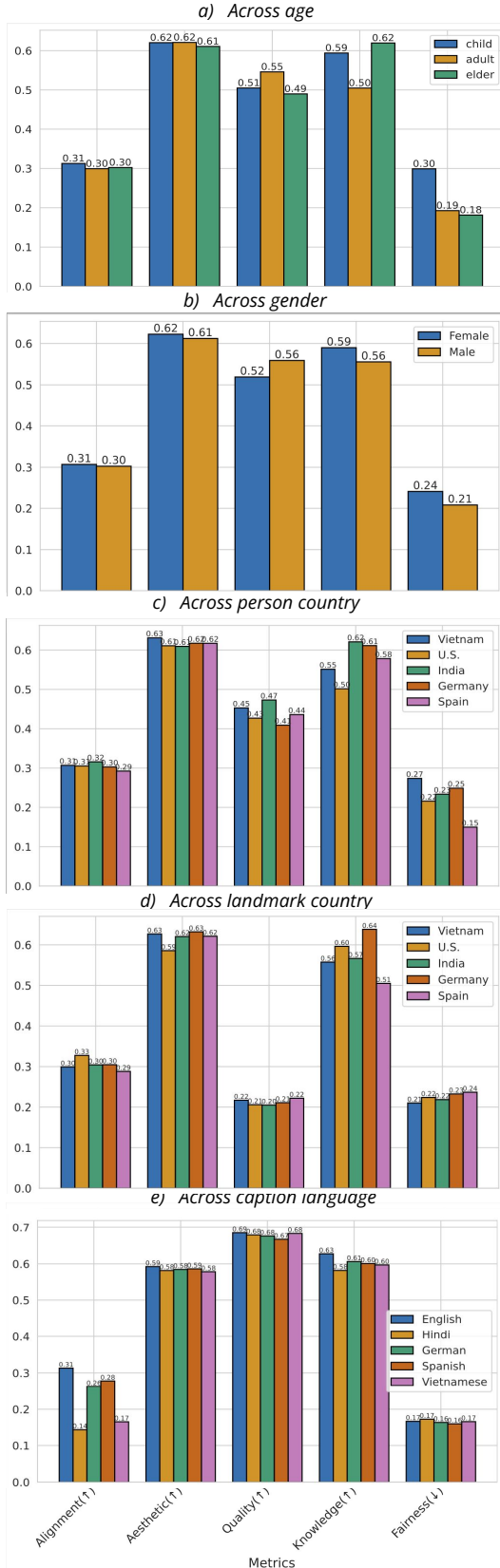


Figure 4: Ablation studies on (a) person age, (b) person gender, (c) person country, (d) landmark country, (e) caption language using the best overall model, the Multi-agent English Flux-M (a-d) and Multi-agent Multilingual Alt-M (e). Performance across all five metrics—Alignment, Aesthetic, Quality, Knowledge, and Fairness—reveals significant variation across these demographic categories.

a) Person Age. Figure 4 a) shows that Image Quality varies by age group, with Adults achieving the highest quality (0.55), followed by Children (0.51) and Elders (0.49). The model is also fairer when depicting Elders (0.18) and Adults (0.19) compared to Children (0.30).

b) Person Gender. Figure 4 b) shows that Knowledge and Image Quality varies by gender, with Males achieving higher quality (0.56) than Females (0.52). However, the model is fairer when depicting Males (0.21) than Females (0.24). The other metrics remain consistent across both groups.

c) Person Country. Figure 4 c) shows that model performance varies by person’s country. Alignment is highest for Indian people (0.32) and lowest for Spanish people (0.29). Similarly, Image Quality is highest for Indian people (0.47) and lowest for German people (0.41). The model is also fairest when depicting Spanish (0.15) and least fair for Vietnamese (0.27).

d) Landmark Country. Figure 4 d) shows that model performance varies by landmark country. The most notable difference is in the Knowledge metric, with German landmarks being the most well-known (0.64), followed by U.S. (0.60), Indian (0.54), Vietnamese (0.50), and Spanish (0.51). Alignment is highest for U.S. landmarks (0.33) and lowest for Spanish landmarks (0.29).

e) Caption Language. Figure 4 e) shows that model performance varies by caption language, with English achieving the highest Alignment (0.31) and Knowledge (0.63), while Hindi and Vietnamese score the lowest (0.14 and 0.43, respectively). This disparity may stem from differences in training data availability, as model performance moderately correlates with dataset size (Pearson coefficient: 0.5), estimated from CommonCrawl (Wenzek et al., 2020). Furthermore, models with English captions achieve higher Alignment than non-English (0.30 vs. 0.20) (see Figure 10).

f) Intersectionality. Examining a single demographic category, such as race or gender, may overlook nuanced inequalities (Field et al., 2021). To address this, we analyze the intersectionality of age and gender, person and landmark country, and language and person country. We measure Alignment and analyze other metrics across various demographic intersections, as detailed in Appendix E.2. **Age and Gender.** Figure 5 (right) shows that Alignment performance varies by gender for generating adult images, with males having a lower

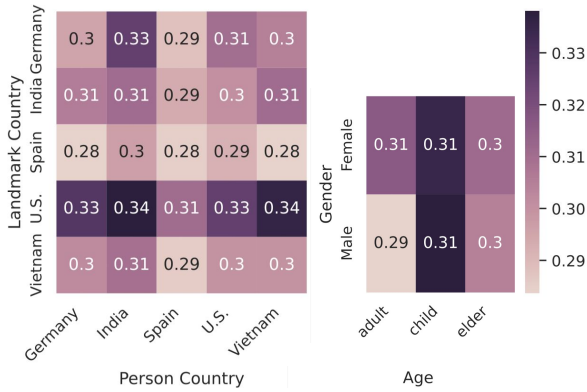


Figure 5: Alignment with best overall model, Flux-M, over person-landmark (left) and gender-age (right).

score (0.29) compared to females (0.31). The performance for child and elder categories remains consistent across gender.

Person and Landmark Country. Figure 5 (left) illustrates Alignment across Person and Landmark Country. We expected higher performance when the person and landmark originate from the same country, suggesting challenges in cross-cultural representation. However, results vary by country. For instance, the highest alignment occurs when Indian or Vietnamese people visit U.S. landmarks (0.34), comparable to U.S. people at U.S. landmarks (0.33). In contrast, the lowest alignment is observed when Vietnamese people visit Spanish landmarks (0.28). All metrics are detailed in Appendix E.2.

Language and Country. Figure 6 shows Alignment across Person Country and Caption Language. English, Spanish, and Vietnamese captions achieve the highest performance (~ 0.3) with minimal variation across person countries. However, Hindi captions perform best for Indian people (0.17) and worst for Spanish and U.S. people (0.13). This suggests that, for certain languages, the interaction between caption language and the depicted person’s culture influences Alignment in image generation.

4.4 Human Evaluation and Error Analysis

Two annotators evaluate a subset of 300 images, covering all demographics (age, gender, country, landmark) and model settings (Alt-S, Alt-M, Flux-S, Flux-M). They assess the generated images based on three key metrics: Alignment, Quality, and Aesthetics. Following Lee et al. (2024), Quality is measured in terms of photorealism, while Aesthetics is evaluated based on subject clarity and overall visual appeal. Annotator agreement is mea-

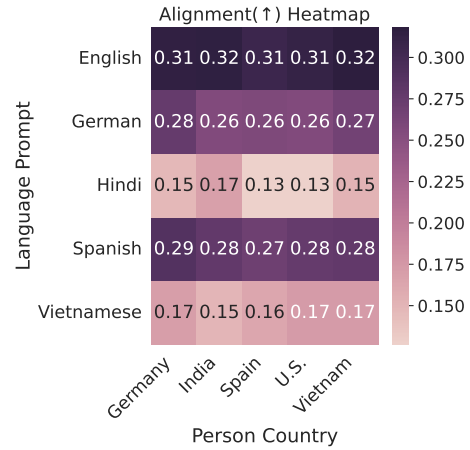


Figure 6: Alignment with best multilingual model, Alt-M, over image caption language and person country.

sured using weighted Cohen’s Kappa for ordinal values (Cohen, 1968), yielding scores between 0.5 and 0.6 across all three metrics, indicating moderate agreement. The complete set of human evaluation questions, along with the annotation interface, is detailed in Appendix D.

Most Common Errors. Across models, errors primarily involve incorrect backgrounds and failures in human rendering. For Flux-M, background inaccuracies are most frequent (38/75), followed by deviations from prompt details and occasional human rendering errors (5/75), such as missing fingers or misplaced cultural markers; landmark errors are comparatively rare (2/75). In contrast, Flux-S exhibits substantially more landmark omissions (15/75) and increased human rendering errors (10/75), particularly for traditional attire and facial features. The Alt models show more severe artifacts overall, with frequent background errors (55/75), pronounced body distortions, and multiplicity errors. While Alt-M reduces culturally related errors (2/75), it still exhibits body distortions (15/75).

4.5 Qualitative Results

In Figure 7, we compare the images generated by our multi-agent framework (Flux-M and Alt-M) with those from simpler models (Flux-S and Alt-S). The second column presents images generated with Vietnamese captions using the multilingual models (Alt-Vi-S, Alt-Vi-M). Compared to the simple models, the multi-agent models perform better at generating landmarks and people. However, they still miss important details about people, such as *a person looking up, curly hair, or hair tied back with a nón lá hat*. Notably, body distortions are more pronounced in the Alt-S model. While the



Figure 7: Comparison of generated images and captions from multi-agent (Flux-M, Alt-M) and simple models (Flux-S, Alt-S). The first two columns show where multi-agent models perform better, while the last column shows where simpler models excel. The second column depicts images generated with *Vietnamese* captions using the multilingual model Alt (Alt-Vi-S, Alt-Vi-M). Demographic keywords are **bolded**, and errors are marked in **red**.

Flux model produces more accurate backgrounds, they tend to be blurrier compared to those in the Alt model. A manual error analysis of 300 images across all demographics highlights the need for further improvements, particularly in rendering body structures and backgrounds. Additional results across demographics are in Appendix E.3.

5 Lessons Learned and Actionable Steps

Our study surfaces insights into the challenges of *multicultural text-to-image generation* and highlights several directions for improving cultural grounding, demographic coverage, and evaluation practices in future models.

Richer Cultural and Demographic Captions.

Multicultural generation benefits from richer prompts that clearly articulate cultural context, demographic attributes, and landmark-specific details. By integrating diverse perspectives through collaboration, multi-agent models enhance alignment, aesthetics, quality, and knowledge (Section 4.2). Future research should focus on refining multi-agent frameworks to further enhance alignment and representational diversity. Our work can also be extended to evaluate a broader range of cultural interactions—such as social activities, rituals, and everyday practices—to better assess reasoning and action-based image generation.

Multilingual Support Remains a Bottleneck.

We observe systematic performance gaps between English and non-English prompts, with English-language generations consistently achieving higher alignment scores (Figure 4e). These disparities point to limitations in current multilingual training and evaluation practices. Improving multilin-

gual coverage—both in training data and model architectures—is essential for achieving equitable performance across languages and cultures.

Develop Better Evaluation Metrics. Automated metrics do not always align with qualitative judgments in multicultural settings, particularly when visually plausible context inflates scores despite incorrect culturally salient elements (e.g., accurate surroundings but an incorrect landmark; Section 4.4). More reliable evaluation requires metrics that place greater emphasis on landmark identity, demographic attributes, and compositional correctness. Targeted automated measures, complemented by human evaluation, remain essential for accurately assessing multicultural image generation.

6 Conclusion

In this paper, we introduce *multicultural text-to-image generation* as a new task for evaluating how models depict people and landmarks from different cultural backgrounds within a single image. We release MOSAIG the first benchmark for this setting, comprising 9,000 images across five countries, three age groups, two genders, 25 historical landmarks, and five languages. Our automated and human evaluations reveal substantial variation across demographics, languages, and cultural configurations, highlighting persistent gaps in multilingual and cross-cultural generation. Overall, our findings emphasize the importance of explicit cultural and demographic grounding for improving image quality, factual correctness, and representational consistency. We release our dataset and evaluation framework to support multicultural text-to-image generation: <https://anonymous.4open.science/r/MosAIG>.

615 Limitations and Ethical Considerations

616 **Scope of Demographic Coverage.** Our study
617 considers a limited set of demographic attributes,
618 focusing on binary gender categories, three coarse
619 age groups (child, adult, elder), and five countries
620 and languages. These choices necessarily simplify
621 the rich diversity of gender identities, life stages,
622 and cultural experiences, and limit our ability to
623 assess performance for underrepresented commu-
624 nities. We view this design as a proof of concept
625 for studying multicultural image generation in a
626 controlled setting. Importantly, our dataset and
627 pipeline are fully open-source and designed to be
628 easily extended to additional countries, languages,
629 age groups, and gender identities in future work.

630 **Challenges in Modeling Cultural Identity.** Our
631 approach relies on structured prompts to approxi-
632 mate cultural and demographic context, but identity
633 is inherently complex and cannot be fully captured
634 through high-level attributes such as nationality,
635 language, age, or gender alone (Saha et al., 2025).
636 Defining culture primarily through national affil-
637 iation risks overlooking substantial intra-cultural
638 variation and lived experience. Future work should
639 incorporate richer contextual dimensions—such as
640 historical background, social practices, and per-
641 sonal narratives—to enable more nuanced and au-
642 thentic representations.

643 **Limitations of Automated Alignment Metrics.**
644 We rely on CLIPScore as a scalable, reference-free
645 measure of text–image alignment, but this metric
646 has several limitations. Its coarse-grained con-
647 trastive training makes it insensitive to fine-grained
648 compositional errors, such as incorrect spatial re-
649 lationships or misattributed attributes, and it may
650 assign high scores to images that contain the correct
651 objects in incorrect configurations (Hessel et al.,
652 2021). CLIPScore is also largely insensitive to
653 word order and linguistic phenomena such as nega-
654 tion, and exhibits biases toward salient or centrally
655 positioned objects, reducing its reliability for com-
656 plex, multi-object prompts (Castro et al., 2023;
657 Abbasi et al., 2025). We therefore encourage fu-
658 ture work to complement CLIPScore with image-
659 based classifiers and targeted evaluation methods
660 that more directly assess visual grounding and de-
661 mographic fidelity (Hu et al., 2023).

References 662

- 663 Reza Abbasi, Ali Nazari, Aminreza Sefid, Moham-
664 madali Banayeezade, Mohammad Hossein Ro-
665 hban, and Mahdieh Soleymani Baghshah. 2025. *An-*
666 *alyzing clip’s performance limitations in multi-object*
667 *scenarios: A controlled high-resolution study*. *ArXiv*,
668 *abs/2502.19828*.
- Muhammad Farid Adilazuarda, Sagnik Mukherjee,
669 Pradhyumna Lavania, Siddhant Shivdutt Singh, Al-
670 ham Fikri Aji, Jacki O’Neill, Ashutosh Modi, and
671 Monojit Choudhury. 2024. *Towards measuring and*
672 *modeling “culture” in LLMs: A survey*. In *Proceed-*
673 *ings of the 2024 Conference on Empirical Methods in*
674 *Natural Language Processing*, pages 15763–15784,
675 Miami, Florida, USA. Association for Computational
676 Linguistics. 677
- Longju Bai, Angana Borah, Oana Ignat, and Rada Mi-
678 halcea. 2025. *The power of many: Multi-agent mul-*
679 *timodal models for cultural image captioning*. In
680 *Proceedings of the 2025 Conference of the Nations*
681 *of the Americas Chapter of the Association for Com-*
682 *putational Linguistics: Human Language Technol-*
683 *ogies (Volume 1: Long Papers)*, pages 2970–2993,
684 Albuquerque, New Mexico. Association for Compu-
685 tational Linguistics. 686
- James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jian-
687 feng Wang, Linjie Li, Long Ouyang, Juntang Zhuang,
688 Joyce Lee, Yufei Guo, and 1 others. 2023. Improving
689 image generation with better captions. *Computer Sci-*
690 *ence*. <https://cdn.openai.com/papers/dall-e-3.pdf>,
691 2(3):8. 692
- Mehar Bhatia, Sahithya Ravi, Aditya Chinchure, Eu-
693 nJeong Hwang, and Vered Shwartz. 2024. *From*
694 *local concepts to universals: Evaluating the multi-*
695 *cultural understanding of vision-language models*.
696 In *Proceedings of the 2024 Conference on Empiri-*
697 *cal Methods in Natural Language Processing*, pages
698 6763–6782, Miami, Florida, USA. Association for
699 Computational Linguistics. 700
- Stephen Castles, Hein de Haas, and Miller Mark J.
701 2103. *The Age of Migration: International Popu-*
702 *lation Movements in the Modern World*. 703
- Santiago Castro, Oana Ignat, and Rada Mihalcea. 2023.
704 *Scalable performance analysis for vision-language*
705 *models*. In *STARSEM*. 706
- Jacob Cohen. 1968. *Weighted kappa: Nominal scale*
707 *agreement provision for scaled disagreement or par-*
708 *tial credit*. *Psychological Bulletin*, 70:213–220. 709
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal,
710 Vishrav Chaudhary, Guillaume Wenzek, Francisco
711 Guzmán, Edouard Grave, Myle Ott, Luke Zettle-
712 moyer, and Veselin Stoyanov. 2020. *Unsupervised*
713 *cross-lingual representation learning at scale*. In *Pro-*
714 *ceedings of the 58th Annual Meeting of the Asso-*
715 *ciation for Computational Linguistics*, pages 8440–
716 8451, Online. Association for Computational Lin-
717 guistics. 718

719	Javier Martín Daniel Verdú. 2024. Flux.1 lite: Distilling flux1.dev for efficient text-to-image generation.	776
720		777
721	Anjalie Field, Su Lin Blodgett, Zeerak Waseem, and Yulia Tsvetkov. 2021. A Survey of Race, Racism, and Anti-Racism in NLP . In <i>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)</i> , pages 1905–1925, Online. Association for Computational Linguistics.	778
722		779
723		780
724		781
725		782
726		783
727		
728		
729	Dhruba Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. 2024. Geneval: An object-focused framework for evaluating text-to-image alignment. <i>Advances in Neural Information Processing Systems</i> , 36.	
730		
731		
732		
733		
734	Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V. Chawla, Olaf Wiest, and Xiangliang Zhang. 2024. Large language model based multi-agents: a survey of progress and challenges . In <i>Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI '24</i> .	
735		
736		
737		
738		
739		
740	Shanshan Han, Qifan Zhang, Yuhang Yao, Weizhao Jin, Zhaozhuo Xu, and Chaoyang He. 2024. Llm multi-agent systems: Challenges and open problems. <i>arXiv preprint arXiv:2402.03578</i> .	
741		
742		
743		
744	Sebastian Hartwig, Dominik Engel, Leon Sick, Hannah Kniesel, Tristan Payer, Timo Ropinski, and 1 others. 2024. Evaluating text to image synthesis: Survey and taxonomy of image quality metrics. <i>arXiv preprint arXiv:2403.11821</i> .	
745		
746		
747		
748		
749	Daniel Hershcovich, Stella Frank, Heather Lent, Miryam de Lhoneux, Mostafa Abdou, Stephanie Brandl, Emanuele Bugliarello, Laura Cabello Piqueras, Ilias Chalkidis, Ruixiang Cui, Constanza Fierro, Katerina Margatina, Phillip Rust, and Anders Søgaard. 2022. Challenges and strategies in cross-cultural NLP . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 6997–7013, Dublin, Ireland. Association for Computational Linguistics.	
750		
751		
752		
753		
754		
755		
756		
757		
758		
759		
760	Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 2021. CLIPScore: A reference-free evaluation metric for image captioning . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 7514–7528, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	
761		
762		
763		
764		
765		
766		
767		
768	Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang, Mari Ostendorf, Ranjay Krishna, and Noah A Smith. 2023. Tifa: Accurate and interpretable text-to-image faithfulness evaluation with question answering. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 20406–20417.	
769		
770		
771		
772		
773		
774	Nithish Kannan, Arif Ahmad, Marco Andreetto, Vinodkumar Prabhakaran, Utsav Prabhu, Adji Bousso	
775		
	Dieng, Pushpak Bhattacharyya, and Shachi Dave. 2024. Beyond aesthetics: cultural competence in text-to-image models. In <i>Proceedings of the 38th International Conference on Neural Information Processing Systems, NIPS '24</i> , Red Hook, NY, USA. Curran Associates Inc.	776
		777
		778
		779
		780
		781
	Black Forest Labs. 2024. Flux. https://github.com/black-forest-labs/flux .	782
		783
	Tony Lee, Michihiro Yasunaga, Chenlin Meng, Yifan Mai, Joon Sung Park, Agrim Gupta, Yunzhi Zhang, Deepak Narayanan, Hannah Teufel, Marco Bellagente, and 1 others. 2024. Holistic evaluation of text-to-image models. <i>Advances in Neural Information Processing Systems</i> , 36.	784
		785
		786
		787
		788
		789
	Zhiqiu Lin, Deepak Pathak, Baiqi Li, Jiayao Li, Xide Xia, Graham Neubig, Pengchuan Zhang, and Deva Ramanan. 2025. Evaluating text-to-visual generation with image-to-text generation. In <i>European Conference on Computer Vision</i> , pages 366–384. Springer.	790
		791
		792
		793
		794
	Bingshuai Liu, Longyue Wang, Chenyang Lyu, Yong Zhang, Jinsong Su, Shuming Shi, and Zhaopeng Tu. 2024. On the cultural gap in text-to-image generation. In <i>ECAI 2024</i> , pages 930–937. IOS Press.	795
		796
		797
		798
	Rada Mihalcea, Oana Ignat, Longju Bai, Angana Borah, Luis Chiruzzo, Zhijing Jin, Claude Kwizera, Joan Nwatu, Soujanya Poria, and Tamar Solorio. 2025. Why ai is weird and shouldn't be this way: towards ai for everyone, with everyone, by everyone . In <i>Proceedings of the Thirty-Ninth AAAI Conference on Artificial Intelligence and Thirty-Seventh Conference on Innovative Applications of Artificial Intelligence and Fifteenth Symposium on Educational Advances in Artificial Intelligence, AAAI'25/IAAI'25/EAAI'25</i> . AAAI Press.	799
		800
		801
		802
		803
		804
		805
		806
		807
		808
		809
	Tarek Naous, Michael J Ryan, Alan Ritter, and Wei Xu. 2024. Having beer after prayer? measuring cultural bias in large language models . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 16366–16393, Bangkok, Thailand. Association for Computational Linguistics.	810
		811
		812
		813
		814
		815
		816
	Tuan-Phong Nguyen, Simon Razniewski, Aparna Varde, and Gerhard Weikum. 2023. Extracting cultural commonsense knowledge at scale. In <i>Proceedings of the ACM Web Conference 2023</i> , pages 1907–1917.	817
		818
		819
		820
	Siddhesh Pawar, Junyeong Park, Jiho Jin, Arnav Arora, Junho Myung, Srishti Yadav, Faiz Ghifari Haznitrana, Inhwa Song, Alice Oh, and Isabelle Augenstein. 2024. Survey of cultural awareness in language models: Text and beyond. <i>arXiv preprint arXiv:2411.00860</i> .	821
		822
		823
		824
		825
		826
	Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. 2023. Sdxl: Improving latent diffusion models for high-resolution image synthesis. <i>arXiv preprint arXiv:2307.01952</i> .	827
		828
		829
		830
		831

Age	Gender	Country	Landmark
Child/ Adult/ Elder	Female/Male	Germany	Cologne Cathedral Reichstag Building Neuschwanstein Castle Brandenburg Gate Holocaust Memorial
		India	Taj Mahal Lotus Temple Gateway of India India Gate Charminar
		Spain	Sagrada Familia Alhambra Guggenheim Museum Roman Theater of Cartagena Royal Palace of Madrid
		U.S.	White House Statue of Liberty Mount Rushmore Golden Gate Bridge Lincoln Memorial
		Vietnam	Meridian Gate of Huế Independence Palace One Pillar Pagoda Ho Chi Minh Mausoleum Thien Mu Pagoda

Table 1: Demographics Overview: 3 Age groups, 2 Genders, 5 Countries, and 25 Landmarks

seed: 11000, and processed 3,750 images of the same resolution in 16 hours, averaging about 15 seconds per image. All processing times accounted for overhead related to model loading and image saving, ensuring consistency in image resolution (768×768 pixels) across both models.

D Human Evaluation and Error Analysis

We rely on human annotators to assess a sample of the generated images based on three key metrics: Alignment, Quality, and Aesthetics. Following Lee et al. (2024), Quality is evaluated in terms of photorealism, while Aesthetics is assessed based on subject clarity and overall visual appeal. The complete set of human evaluation questions is outlined below. Annotators are provided with definitions (Table 2) and corresponding questions to guide their assessments. To determine whether the generated images meet their expectations, we ask annotators to rate them using a 5-point Likert scale.

Alignment. We ask the annotators to rate how well the image matches the description.

How well does the image match the description?

1. Does not match at all
2. Has significant discrepancies
3. Has several minor discrepancies
4. Has a few minor discrepancies
5. Matches exactly

Quality. We ask the annotators to rate how photorealistic the generated images are.

Determine if the following image is AI-generated or real.

1. AI-generated photo.
2. Probably an AI-generated photo, but photorealistic.
3. Neutral.
4. Probably a real photo, but with irregular textures and shapes.

Conv. Round	Agent Role	Prompt
Round 1	Country Agent	SYSTEM: You are a {nationality} person from {country} who knows the culture of this country well. USER: Provide a visual description of culturally appropriate traditional clothing, accessories, and colors, for the {nationality} person. Focus on specific materials, key cultural patterns, and symbolic colors. Your response must be under 25 words. \nASSISTANT:
	Landmark Agent	SYSTEM: You are a person who has visited {place} many times and know this landmark well. USER: Provide a visual description of its architectural features, colors, and environmental details. Your response must be under 25 words. \nASSISTANT:
	Age-Gender Agent	SYSTEM: You are a {age_gender_combined} and can describe traits of this person well. USER: Provide a visual description of attire, accessories, and physical details. Focus on skin, body, hair texture, and accessories. Your response must be under 25 words. \nASSISTANT:
Round 2	Country Agent	SYSTEM: You are a {nationality} person from {country}. USER: Enhance the persona description by addressing: 'How would a person's clothing harmonize with the colors of {place}?'. Ensure cultural significance is highlighted. \nASSISTANT:
	Landmark Agent	SYSTEM: You are a person who knows {place} well. USER: Enhance the place description by addressing: 'What visual elements of {place} would complement the persona's attire?'. Limit to under 25 words. \nASSISTANT:
	Age-Gender Agent	SYSTEM: You are a {age_gender_combined}. USER: Enhance the age-gender description by addressing: 'What attire adjustments could reflect age-appropriate traits for a {nationality} {age_gender_combined}?'. Ensure specific details on attire and physical traits. \nASSISTANT:
Round 3	Summarizer Agent	SYSTEM: You excel at crafting concise visual prompts. USER: Give a final prompt in a single line under 48 words and under 77 tokens strictly. Ensure the words {nationality} and {age_gender_combined} of the person and other descriptions with the {place} background are mentioned explicitly in the final prompt. \nASSISTANT:

Figure 8: Our Multi-agent Framework Prompts

Aspect	Definition
Alignment	Is the image semantically correct given the text (text-image alignment)?
Quality	Do the generated images look like real photographs?
Aesthetic	Is the image aesthetically pleasing?
Fairness	Does the model exhibit performance disparities across social groups (e.g., gender, dialect)
Knowledge	Does the model have knowledge about the world or domains?

Table 2: Evaluation Aspects of Text-to-Image Models

5. Real photo.

5. The image is aesthetically stunning. I can look at it all day.

Aesthetics. To evaluate the overall aesthetics, we ask annotators to provide a holistic assessment of the image’s visual appeal by rating its aesthetic quality.

How aesthetically pleasing is the image?

1. I find the image ugly.
2. The image has a lot of flaws, but it’s not completely unappealing.
3. I find the image neither ugly nor aesthetically pleasing.
4. The image is aesthetically pleasing and is nice to look at.

990
991

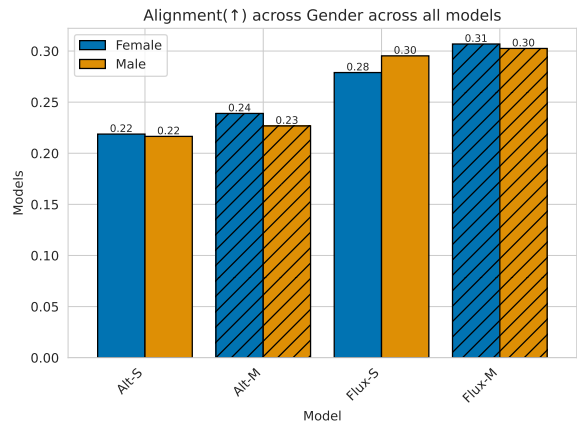
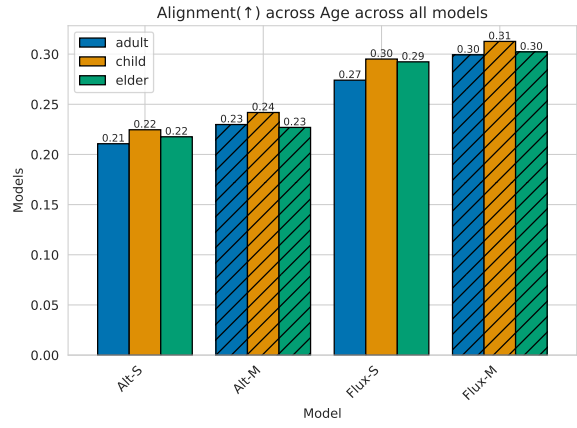


Figure 9: Human Annotation Interface for manually evaluating the images across all models.

992
993
994

E Results

E.1 Across Metrics and Demographics, across All Models

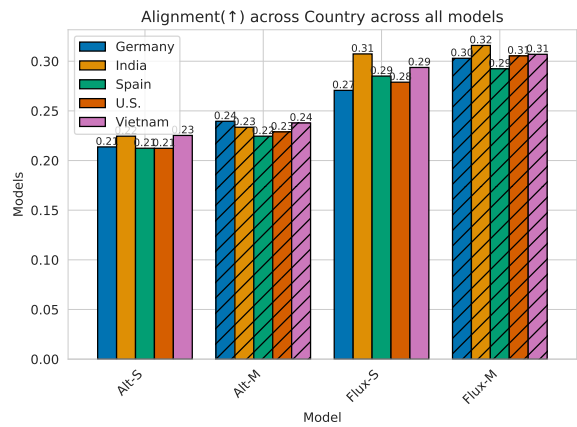


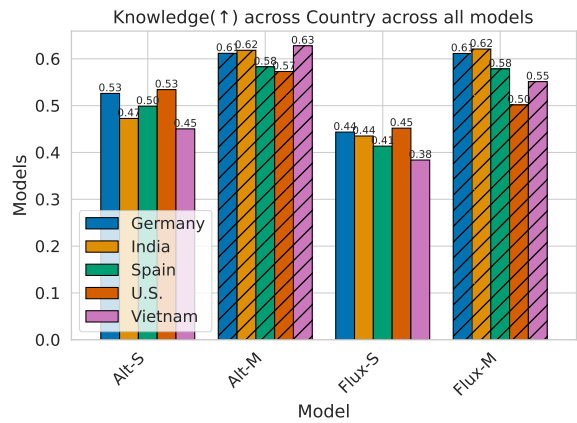
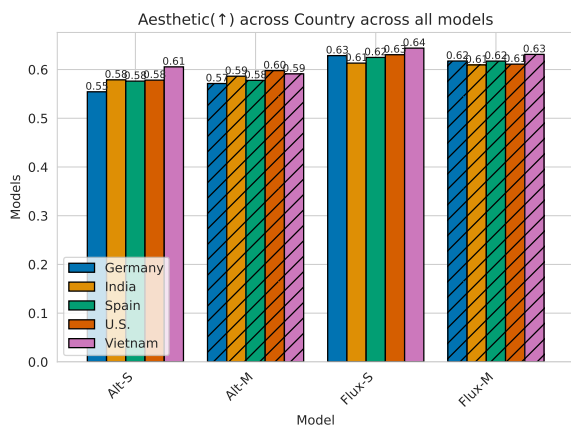
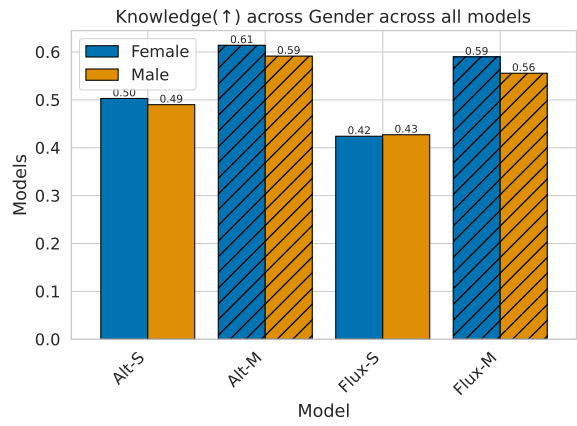
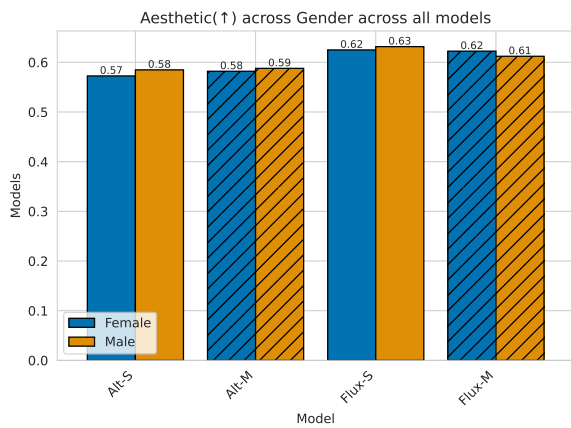
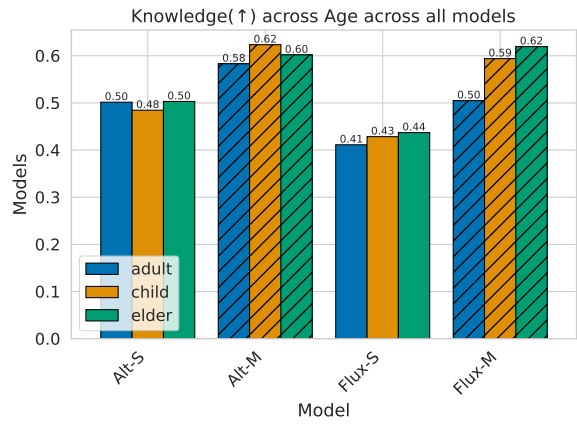
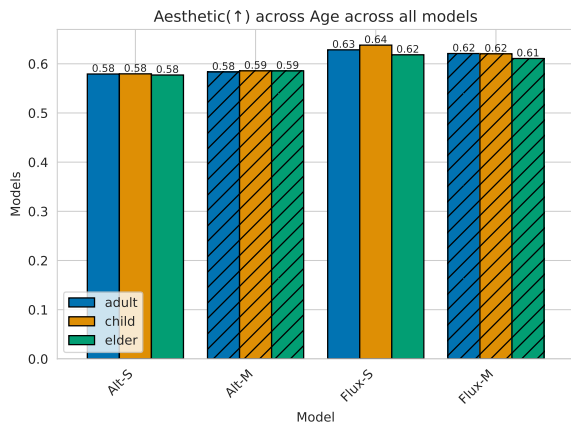
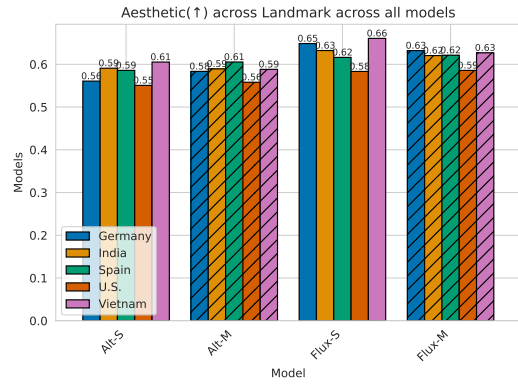
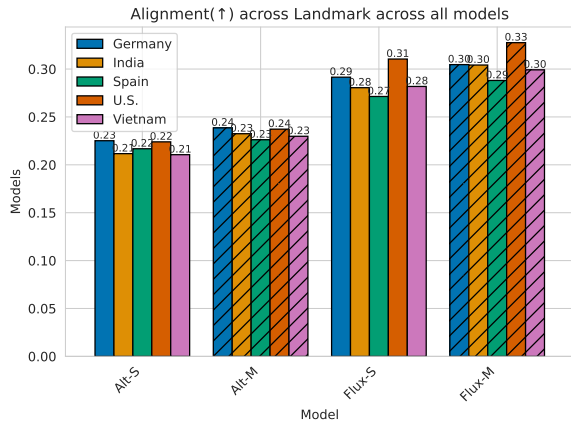
E.2 Intersectionality

995

E.3 Qualitative Results

996





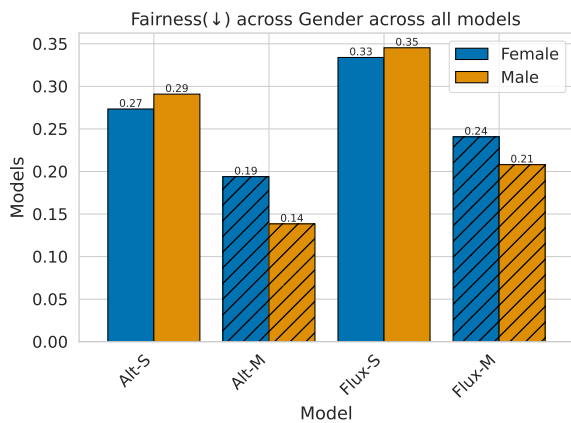
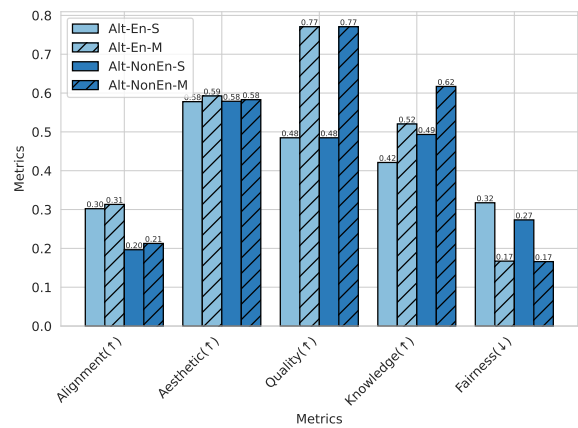
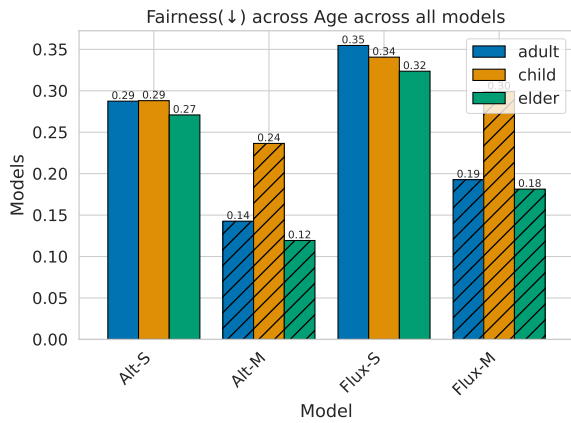
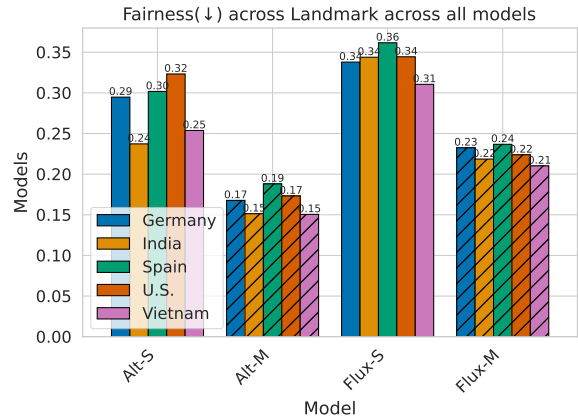
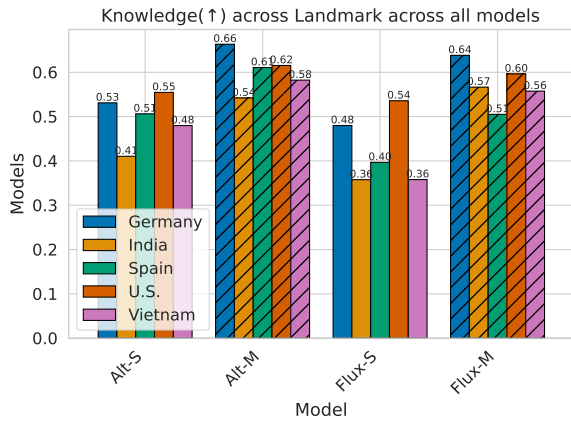
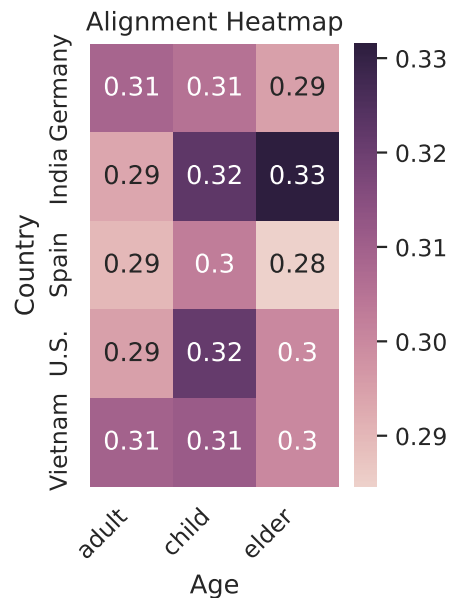
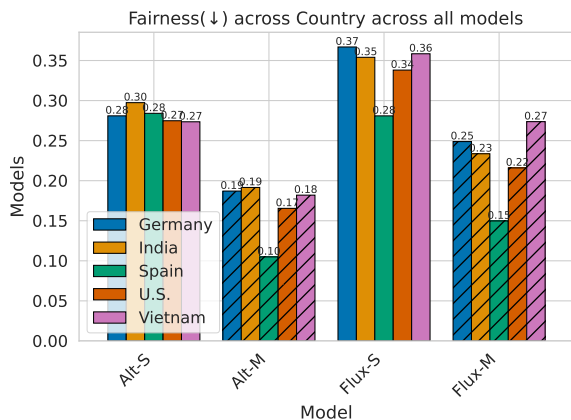
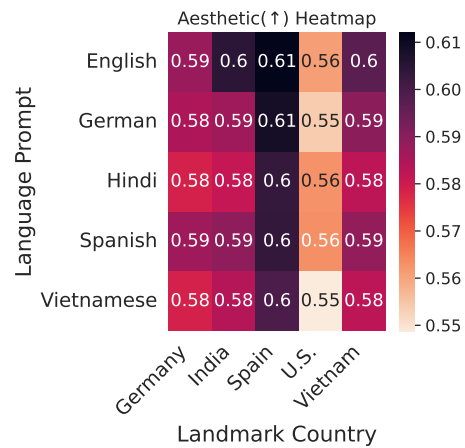
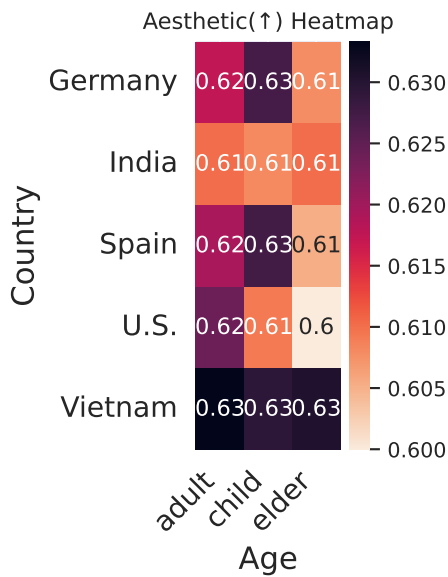
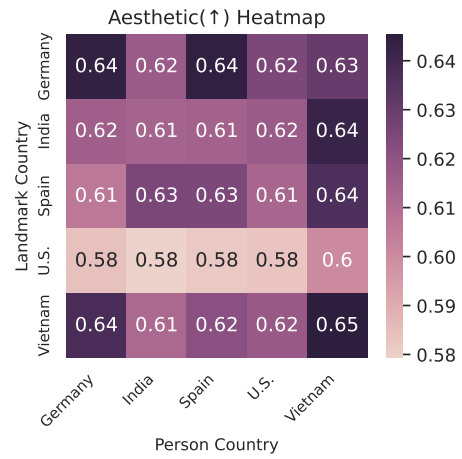
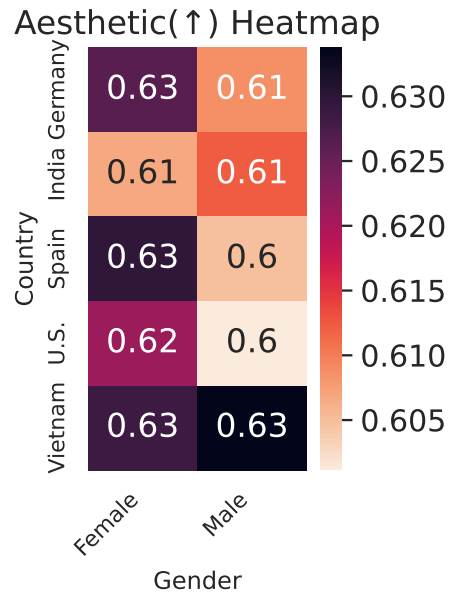
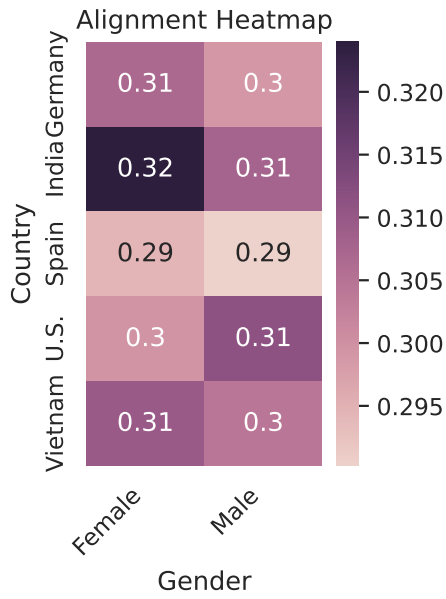
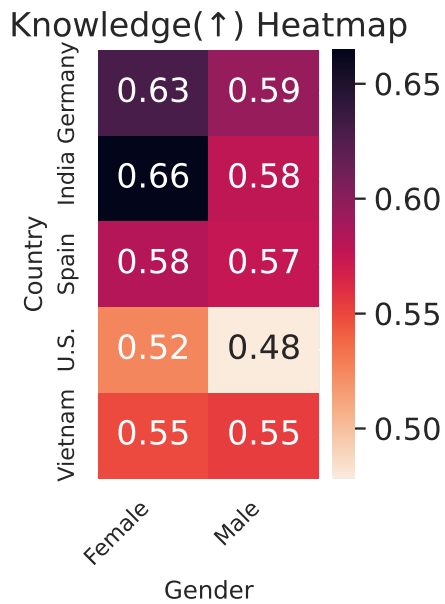
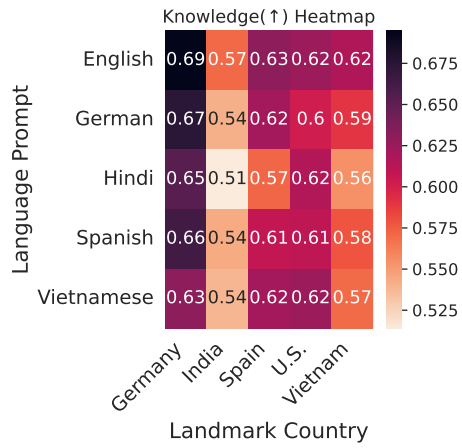
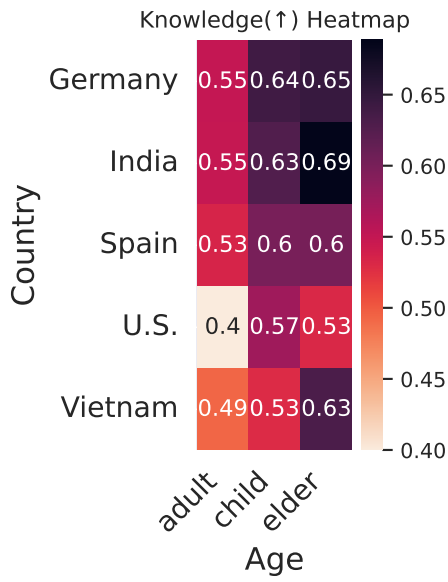
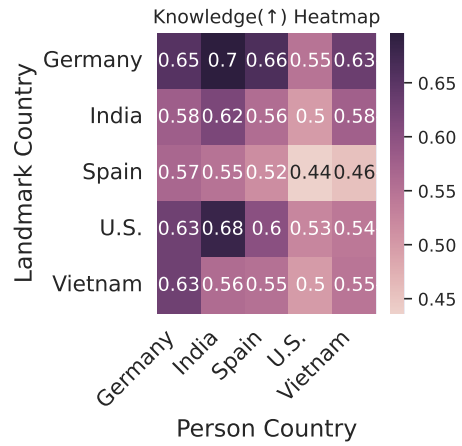
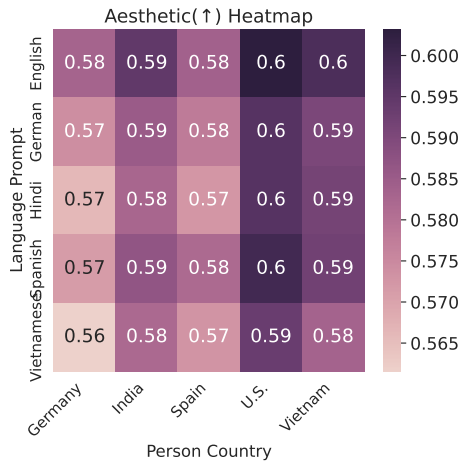


Figure 10: English vs. Multilingual Performance. Models with English captions as input (Alt-En-S, Alt-En-M) achieve higher scores than non-English (Alt-NonEn-S, Alt-NonEn-M) in Alignment (0.30 vs. 0.20), while performing comparably across Aesthetics and Quality metrics. Knowledge and Fairness performance is higher for non-English models







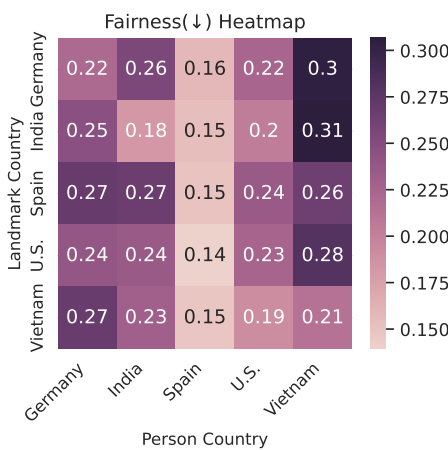
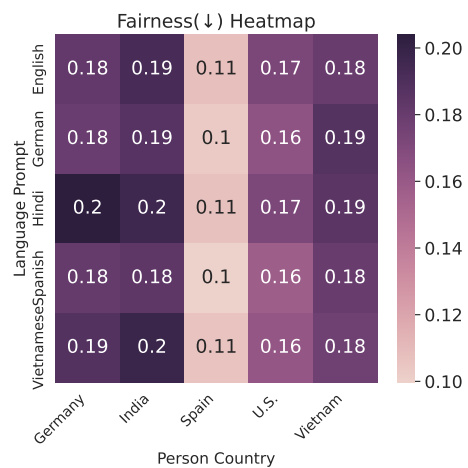
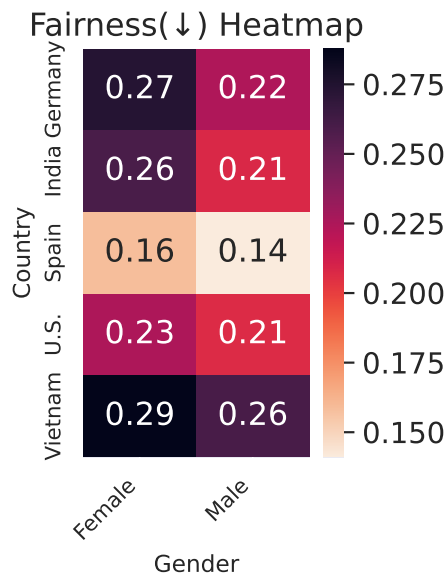
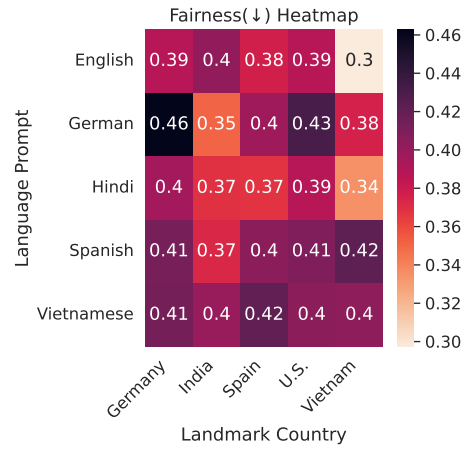
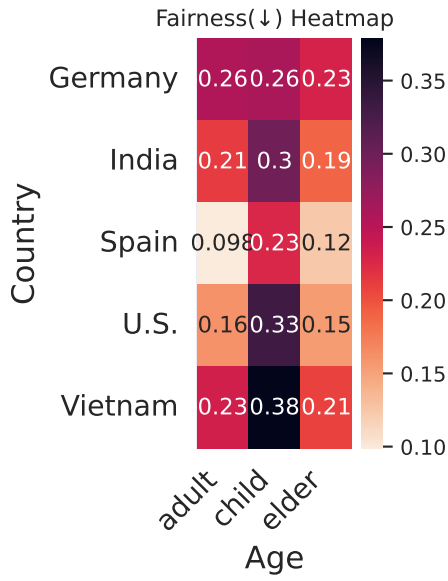




Figure 11: Comparison of generated images and captions using our multi-agent framework (Flux-M, Alt-M) and simple models (Flux-S, Alt-S). The second column depicts images generated with **German** captions using the multilingual model Alt (Alt-De-S, Alt-De-M). Demographic keywords are **bolded**, and incorrect content is marked in **red**.



Figure 12: Comparison of generated images and captions using our multi-agent framework (Flux-M, Alt-M) and simple models (Flux-S, Alt-S). The second column depicts images generated with **Hindi** captions using the multilingual model Alt (Alt-Hi-S, Alt-Hi-M). Demographic keywords are **bolded**, and incorrect content is marked in **red**.



Figure 13: Comparison of generated images and captions using our multi-agent framework (Flux-M, Alt-M) and simple models (Flux-S, Alt-S). The first column depicts images generated with **German** captions using the multilingual model Alt (Alt-De-S, Alt-De-M). The last column depicts images generated with **Spanish** captions using the multilingual model Alt (Alt-Es-S, Alt-Es-M). Demographic keywords are **bolded**, and incorrect content is marked in **red**.