

CULTURALFRAMES: Assessing Cultural Expectation Alignment in Text-to-Image Models and Evaluation Metrics

Anonymous Authors¹

Abstract

The increasing ubiquity of text-to-image (T2I) models as tools for visual content generation raises concerns about their ability to accurately represent diverse cultural contexts. In this work, we present the first study to systematically quantify the alignment of T2I models and evaluation metrics with respect to both *explicit as well as implicit cultural expectations*. To this end, we introduce CULTURALFRAMES, a novel benchmark designed for rigorous human evaluation of cultural representation in visual generations. Spanning 10 countries and 5 socio-cultural domains, CULTURALFRAMES comprises 983 prompts, 3,637 corresponding images generated by 4 state-of-the-art T2I models, and over 10k detailed human annotations. We found that state-of-the-art T2I models not only fail to meet the implicit expectations which are more challenging to meet, but also the less challenging explicit expectations. Across models and countries, cultural expectations are missed an average of 44% of the time. Among these failures, explicit expectations are missed at a surprisingly high average rate of 68%, while implicit expectation failures are also significant, averaging 49%. Furthermore, we demonstrate that existing T2I evaluation metrics correlate poorly with human judgments of cultural alignment, irrespective of their internal reasoning. Collectively, our findings expose critical gaps, providing actionable directions for developing more culturally informed T2I models and evaluation methods.

1. Introduction

Visual media such as advertisements, posters, and public imagery play a central role in encoding and transmitting

cultural values (McLuhan, 1966). They often depict culturally specific elements (e.g., traditional attire, religious symbols) and embed societal norms and values (e.g., expectations around family structure, gender roles, and etiquette), thus reflecting and influencing the cultures from which they originate (Hall, 1980).

Text-to-image (T2I) models are emerging as a significant component of this visual media ecosystem, now adopted across diverse domains like education, marketing, and storytelling (Dehouche & Dehouche, 2023; Loukili et al., 2025; Maharana et al., 2022). This magnifies the cultural implications of their outputs for global audiences (Wan et al., 2024; Hartmann et al., 2025) and raises a critical question: how accurately, and with what depth, do these models depict diverse cultures? While T2I models may generate visually plausible outputs for cultural prompts (e.g., “a bride and groom exchanging vows at their Hindu wedding,” Fig. 1), they often capture explicit details at the expense of crucial, implicit elements integral to the cultural context (such as a sacred fire or officiating priest). Indeed, T2I model performance hinges on accurate cultural representation, which can foster familiarity and trust. Inaccuracies, however, risk reinforcing stereotypes, exclusion, or propagating dominant narratives (Naik & Nushi, 2023).

This necessitates evaluation practices that not only verify faithfulness to the explicit expectations (expectations based on the words in the prompt) but also assess the inference and contextualization of implicit cultural expectations (expectations based on the cultural context mentioned in the prompt). However, current T2I evaluation methodologies predominantly focus on the former by assessing explicit prompt-image consistency using automated metrics (Hu et al., 2023; Hessel et al., 2021; Ku et al., 2024a).¹ Further, existing benchmarks for evaluating T2I models are designed around prompts that emphasize attributes like realism (Saharia et al., 2022), compositionality (Huang et al., 2023; 2025), and safety (Lee et al., 2023), typically using generic or Western-centric prompts. Consequently, current evaluation methods and benchmarks lack adequate representation

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the ICML 2025 Workshop on Models of Human Feedback for AI Alignment. Do not distribute.

¹The only prior work evaluating appropriate contextualization of sensitive content is Akbulut et al. (2025), which focuses on image-to-text for historical events.

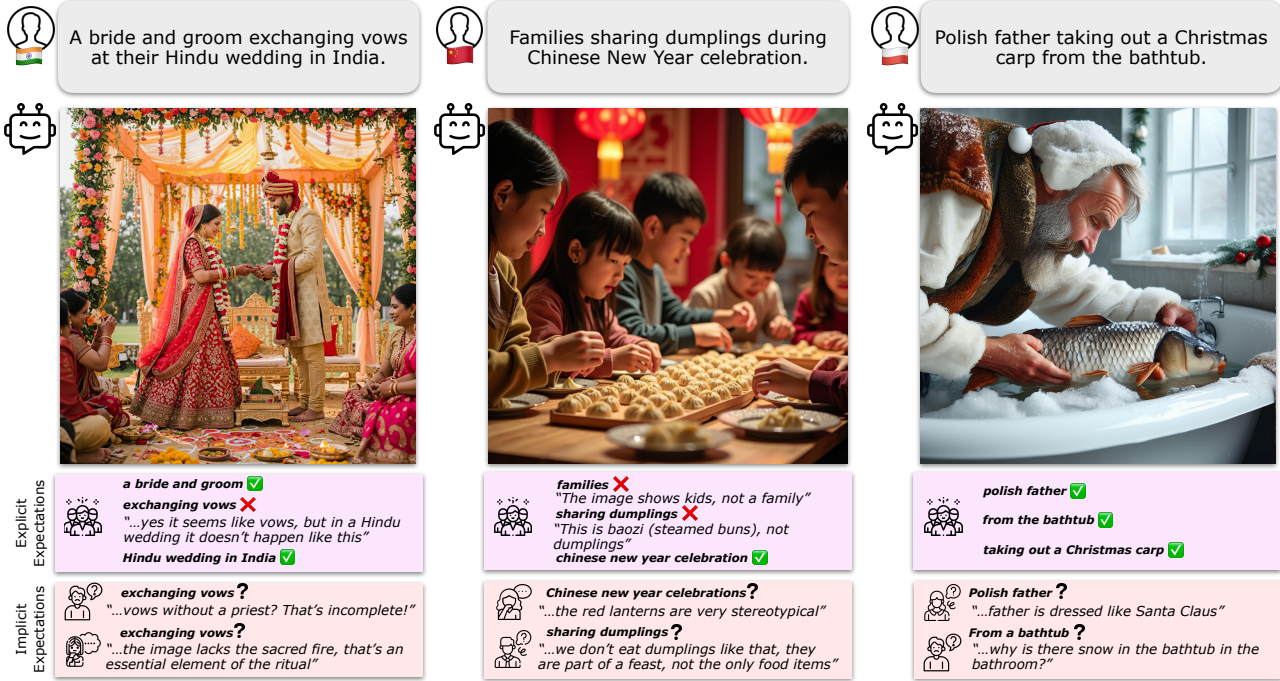


Figure 1: Examples from CULTURALFRAMES benchmark for three selected countries: India, China, and Poland. We ask annotators to evaluate the generated images with respect to both explicit and implicit cultural expectations.

of culturally nuanced and expectation-rich scenarios critical to diverse cultural contexts.

In response to these limitations, we perform a comprehensive study to evaluate how state-of-the-art T2I models represent cultural expectations across diverse contexts. We introduce CULTURALFRAMES, a novel benchmark comprising 983 prompts across 10 countries, with 3,637 corresponding images generated by 4 state-of-the-art T2I models, and over 10k detailed human annotations. The curated prompts are grounded in real-life situations and cover five culturally significant domains: *greetings*, *etiquette*, *dates of significance*, *religion*, and *family life*, which are explicitly designed to test representation of both *explicit* and *implicit* cultural expectations. Using the collected prompts, we first generate images with four state-of-the-art T2I models, two open-source and two closed-source. Second, we conduct evaluations employing human evaluators with relevant cultural backgrounds, who provide fine-grained judgments of the generated images with respect to the prompt in order to assess T2I models' performance. We find that state-of-the-art T2I models not only fail to meet the implicit expectations that are more challenging to meet, but also the less challenging explicit expectations. In fact, models fail to meet cultural expectations 44% of the time on average across countries. Among these instances, the failure rate for explicit expectations is unexpectedly high, averaging 68%, and the rate for implicit expectations is also significant at an average of 49%.

Furthermore, we correlate these human assessments with existing T2I evaluation metrics to demonstrate that current metrics correlate poorly with human judgments of cultural alignment, while differing in their internal reasoning. Collectively, our findings lead to a discussion on actionable directions for developing more culturally informed T2I models and evaluation methodologies. These include utilizing our prompts for future evaluations, leveraging the full CULTURALFRAMES (prompts, images, and annotations) for model alignment, and using explicit instructions for metrics.

2. Related Work

Evaluating T2I models. A suite of benchmarks has been proposed for text-to-image generation. DrawBench (Saharia et al., 2022) and PartiPrompts (Yu et al., 2022) evaluate overall image fidelity and complex scene rendering. The T2I-CompBench series (Huang et al., 2023; 2025) focus specifically on compositional challenges. Human assessment and considerations for bias and fairness are addressed by ImagenHub (Ku et al., 2024c), HEIM (Lee et al., 2023), and GenAI Arena (Jiang et al., 2024). Traditional metrics assess image quality and diversity using embedding-based metrics, e.g., FID (Heusel et al., 2018), Inception Score (Salimans et al., 2016), and the text-image alignment via pre-trained vision-language embeddings, e.g., CLIPScore (Hes-sel et al., 2021) and DinoScore (Ruiz et al., 2023). More

recently, reward models trained on human preferences such as HPSv2 (Wu et al., 2023a), ImageReward (Xu et al., 2023), and PickScore (Kirstain et al., 2023) have shown improved correlation with human judgments. Concurrently, further metrics leverage LLMs and VLMs for evaluating prompt consistency and image-text alignment through question-answering or reasoning, such as TIFA (Hu et al., 2023), DSG (Cho et al., 2024), V2QA (Yarom et al., 2023), VQAScore (Lin et al., 2025), VIEScore (Ku et al., 2024b), and LLMscore (Lu et al., 2023).

Cultural Alignment Evaluation of T2I models. T2I models struggle to accurately and respectfully represent cultural elements, leading to misrepresentation of culturally grounded concepts and values (Ventura et al., 2024; Prabhakaran et al., 2022; Struppek et al., 2023). A growing body of work highlights various cultural biases, such as nationality-based stereotypes (Jha et al., 2024), skin tone bias (Cho et al., 2023), broader risks and social biases in T2I models across gender, race, age, and geography (Bird et al., 2023; Naik & Nushi, 2023). Other works focus on geographic representation (Basu et al., 2023; Hall et al., 2024), showing skewed generations towards Western contexts. Several recent benchmarks aim to probe cultural alignment in T2I systems. CUBE (Kannen et al., 2025) evaluates generations across food, clothing, and landmarks from eight countries. CULTDIFF (Bayramli et al., 2025) studies culturally specific generations across ten nations. CCUB (Liu et al., 2024) introduces a benchmark for inclusive representation and proposes the SCoFT method to leverage model biases for improved equity. Similarly, MC-SIGNS (Yerukola et al., 2025) presents a dataset of gestures from 85 countries, while tasks like cultural image transcreation (Khanuja et al., 2024), study cultural adaptation, evaluating how well models translate images across cultures. Other works retrieve cultural context to refine generation prompts (Jeong et al., 2025), or evaluate portrayals of nationality in limited settings (Alsudais, 2025).

While these efforts provide valuable insights, they predominantly focus on visible and explicit cultural symbols and references like clothing, food, or monuments. Our work is inspired by Qadri et al. (2025), who argue that relying predominantly on standard metrics of faithfulness and quality can yield only surface-level understanding. Therefore, Qadri et al. (2025) advocate for “thick” evaluations, offering qualitative insights through culturally grounded human studies. As a result, our work targets day-to-day scenarios and investigates how well T2I models represent both explicit and implicit cultural expectations. We also evaluate both models and metrics through detailed human studies to understand their strengths and limitations in these scenarios. To the best of our knowledge, this is the first attempt to systematically quantify the alignment of T2I models and metrics with

implicit cultural expectations in visual generations.

3. CULTURALFRAMES

We detail our entire data collection pipeline below and highlight the design decisions that make it distinct from standard annotation efforts.

3.1. Selection of Countries

We operationalize cultural groups using countries as a proxy (Adilazuarda et al., 2024), building upon the premise that individuals within a country share a substantial amount of common cultural knowledge, implicit understandings, and norms that shape their daily interactions and practices (Hofstede et al., 2010; Hershcovich et al., 2022). To create a dataset with diverse cultures, we selected countries spanning five continents and representing diverse cultural zones as per the zone categorization in the World Values Survey (WVS; Haerpfer et al. 2022). Thus, our selection includes countries from the following cultural zones: West and South Asia (India), Confucian (China, Japan), African-Islamic regions (Iran, South Africa), Latin America (Brazil, Chile), English-speaking (Canada), Catholic Europe (Poland), and Protestant Europe (Germany).²

3.2. Selection of Cultural Categories

Our dataset is designed to evaluate culturally relevant expectations in visual generations. Specifically, we target five socio-cultural domains from CulturalAtlas (Mosaica, 2024) deeply embedded in day-to-day life: 1) family, addressing familial roles, hierarchy, and interactions; 2) greetings, covering norms in social and business interactions; 3) etiquette, involving conduct during visits, meals, gift-giving, etc.; 4) religion, reflecting rituals and customs shaping group identities; 5) and dates of significance, highlighting celebrations of cultural, historical, or religious importance. These categories were selected due to their coverage in the CulturalAtlas for the selected countries and their potential to induce prompts that elicit both explicit (elements directly mentioned in the prompt) and implicit cultural (not mentioned in the prompt but inferred from shared cultural commonsense and needed for cultural authenticity) expectations.

3.3. Data Generation Pipeline

Building on cultural categories, we first generate culturally grounded prompts reflecting the core values described above. For each prompt, we generate corresponding images and evaluate across multiple dimensions from culturally knowledgeable annotators to assess whether text-to-image models

²We acknowledge that the labels assigned to these cultural categories are limited in their precision. Yet, these categories present the cross-cultural variation relevant to this work.

Assertion (CulturalAtlas)	Generated Prompts
Greetings (India): Indians expect people to greet the eldest or most senior person first. When greeting elders, some may touch the ground or the elder's feet as a sign of respect.	(1) Grandchildren touching grandfather's feet at an Indian temple. (2) Indian village elder blessing children during harvest festival.
Religion (Iran): Most Iranians believe in Islam, but due to politicization, many younger citizens have withdrawn. Devout followers often practice privately at home.	(1) Iranian family praying together at home. (2) Elderly Iranian man praying in a quiet mosque.

Table 1: Examples of assertions in CulturalAtlas for two categories greetings in India and religion in Iran and corresponding generated prompts.

capture both explicit and implicit cultural expectations.

Prompt Generation. We use Cultural Atlas (Mosaica, 2024) as our knowledge base to extract cultural expectations (norms, practices, values) written as assertions. Cultural Atlas is an educational resource informed by extensive community interviews and validated by cultural experts. To generate culturally grounded prompts, we first extract concise assertions from Cultural Atlas content and feed them to GPT-4o (OpenAI, 2024) using designed instructions (see App. A.1.1). These instructions guide the model to embed cultural expectations into the prompts for realistic and observable everyday scenarios. Next, we use GPT-4o (OpenAI, 2024) and Gemini (Team, 2024) to automatically validate the generated prompts, discarding any that are overly abstract, culturally misaligned, or not visually depictable. As a final step, we present each prompt to three culturally knowledgeable annotators. Only prompts agreed upon by the majority are retained in the dataset (more details in App. A.1.2). Example assertions and prompts from our benchmark are shown in Tab. 1.

Image Generation. We generate images using four state-of-the-art text-to-image models: two open-source models (Flux 1.0-dev (Labs, 2024) and Stable Diffusion 3.5 Large (SD) (Esser et al., 2024)) and two closed-source models (Imagen3 (Imagen-Team-Google, 2024) and GPT-Image (OpenAI, 2025)). We note that Imagen3 includes a prompt expansion mechanism, which we enable by default and also ablate by disabling it to assess its effect on the depiction of cultural expectations. Not focusing on output diversity, we generate one image per model per prompt to keep the evaluation practical. In Fig. 9, we present prompt-image examples.

Rating Collection. We developed a human rating collection interface and the associated annotation guidelines. We tested several interface designs and variants of annotation guidelines to collect high-quality annotations. The final interface and the guidelines are provided in App. A.2. To ensure high data quality, we filtered for attentive annotators and ensured a minimum of 20 unique, culturally knowl-

edgeable workers³ per country. We collect data from three annotators for each country using the Prolific⁴ platform. Our annotation process captures detailed, multi-faceted feedback. Each annotator first evaluates how well the image aligns with the prompt (image-prompt alignment), considering both explicit elements stated in the prompt and implicit elements expected based on cultural context. Following Ku et al. (2024c), we use a 3-point Likert scale: 0.0 (no alignment), 0.5 (partial), and 1.0 (complete). For scores below 1, annotators specify whether explicit, implicit, or both types of elements were missing or not depicted satisfactorily in the image, and highlight the specific words in the prompt whose visual depictions were not satisfactory, along with providing justifications for why they were not satisfactory. This fine-grained rating scheme allows us to analyze the interplay between various quality aspects and their correlation with perceived cultural appropriateness. Annotators flag stereotypes in the images, providing justifications if present. Next, they assess image quality, noting issues such as distortions, artifacts, or unrealistic object rendering. Finally, they assign an overall image score on a 5-point Likert scale.

4. Data Analysis

Prompts. CULTURALFRAMES consists of 983 prompts collected from 10 countries, with each country contributing between 90 and 110 prompts, ensuring balanced cross-country representation. The prompts are distributed across five cultural categories introduced in § 3.2: etiquette (24.3%), religion (14.4%), family (14.2%), greetings (13.1%), and dates of significance (34%). For a detailed per-country breakdown, see Fig. 8 in App. A.1.3.

Images. For open-source models, we generate images for all prompts. However, closed-source models apply safety filters that block some generations. This issue is most noticeable with Imagen3, which filters out 290 prompts—29.5% of the prompts. Most of these are blocked because the prompts involve children. We requested an exemption but have not received approval yet. We will continue to follow up and add more images if access is granted. GPT-4o blocks only 5 prompts. In total, we collect 3,637 images.

Inter-rater Agreement. We collect a total of 10,911 ratings, with each image rated by 3 annotators. To measure agreement among raters, we compute Krippendorff's alpha (Krippendorff, 2013): 0.37 for prompt alignment, 0.28 for image quality, and 0.36 for overall score. These values in-

³Annotators were selected based on the following criteria: born in the country, national of the country, have spent the majority of the first 18 years of life there, and are a resident of the country. The residency criterion was relaxed for China to ensure a sufficient annotator pool size.

⁴<https://www.prolific.com/>

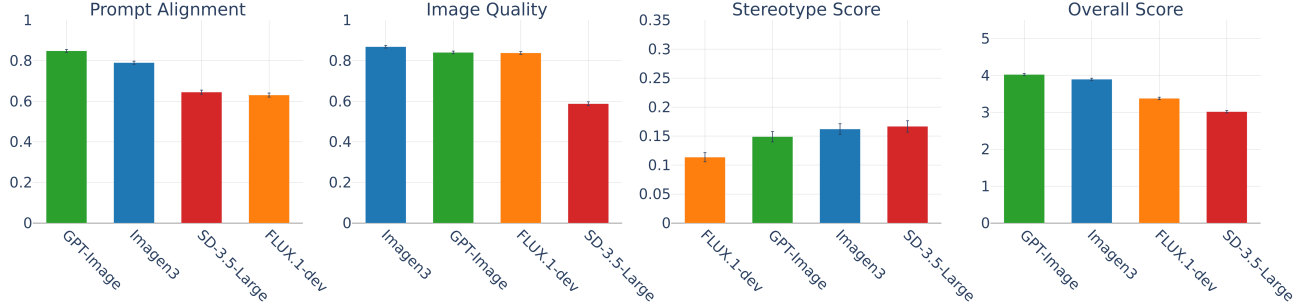


Figure 2: Human evaluation results for selected T2I models. From left to right: 1) Prompt Alignment (0 – 1 scale, 1 =perfect alignment). 2) Image Quality (0 – 1 scale, 1 =highest quality). 3) Stereotype Score (0 – 1 scale, 0 indicates no stereotyping). 4) Overall Score (1 – 5 Likert scale, 5 =best overall). For fairness, we compare across prompts that have images generated by all models.

indicate moderate agreement among annotators. Our results align with previous findings that image quality assessment is subjective (Wu et al., 2023b; Qadri et al., 2025). For prompt alignment, the agreement scores indicate diverse annotators’ expectations, showing the difficulty of the cultural expectation evaluation task.

What aspect of the generated image dominates annotators’ overall assessment? We find that the overall score given by annotators is strongly correlated with image–prompt alignment (Spearman rank correlation of 0.68), whereas image quality shows a more moderate correlation of 0.45. This trend holds consistently across countries, suggesting that annotators prioritize faithfulness to the prompt over aesthetic appeal when rating images. Also, stereotype is negatively correlated with overall score weakly (-0.21), which indicates a lower impact of the presence of stereotypes on overall score. Interestingly, the results contrast with findings from prior work using side-by-side image comparisons (Kirstain et al., 2023), where image quality often dominates overall preference judgments.

5. Evaluating T2I Models on CULTURALFRAMES

How do different models perform for different criteria across different countries? Fig. 2 shows human evaluation results for prompt alignment, image quality, stereotype, and overall score. We find that GPT-Image achieves the highest prompt alignment (0.85), followed by Imagen3 (0.79). The open-source models, SD-3.5-Large and Flux, fall behind with scores of 0.66 and 0.63, respectively. For image quality, Imagen3 is rated highest, with GPT-Image and Flux performing comparably well. SD-3.5-Large, however, scores far behind the other models. Across all models, including the state-of-the-art closed-source ones, the proportion of images rated stereotypical ranged from 10% to

16%, with SD-3.5-Large generating stereotypical visuals the most and Flux the least. Overall, raters prefer images from GPT-Image, consistent with the prompt alignment result. SD received the lowest overall score, most likely due to poorer image quality and higher stereotype levels, despite outperforming Flux in prompt alignment.

Consistent with Rastogi et al. (2024), our findings (Fig. 14) indicate notable cross-country variations in both the overall score and perceived importance of different evaluation criteria. For instance, even assessments of image quality differ, showing a discernible trend where Asian countries tend to assign lower scores across multiple criteria.

Is there a preferred model across countries? For prompt alignment (see Fig. 3), GPT-Image is consistently preferred across countries, followed by Imagen3. Among open-source models, SD-3.5-Large is generally more faithful except for Germany, Poland, and Iran, where Flux performs better. In Fig. 14, we show detailed results across countries and all categories. Regarding image quality, Imagen3 is the preferred model, likely due to its hyper-realistic generations. Interestingly, concerning stereotypes, closed-source models are ranked as more stereotypical for 6 out of the 10 countries.

Which aspect—implicit or explicit—do models fail to capture, and is this consistent across countries? Across CULTURALFRAMES, annotators gave sub-perfect scores (below 1) for 44% of the time. Out of these, 50.3% are attributed to issues with explicit elements, 31.2% to implicit elements, and 17.9% to both. While explicit errors are most common, implicit cultural failures still account for 49.1% of these cases, underscoring persistent challenges in capturing culturally nuanced, context-dependent knowledge. Fig. 4 shows that GPT-Image has the lowest overall image-prompt alignment error rate (ratings $\neq 1$), with its errors roughly evenly split between implicit and explicit types. In contrast, other models, particularly SD-3.5-Large and FLUX, exhibit

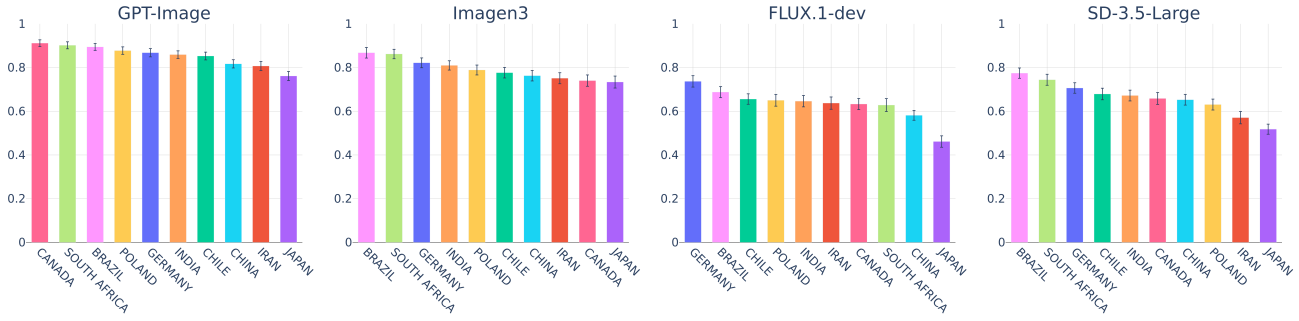


Figure 3: Prompt alignment scores across countries for a given model

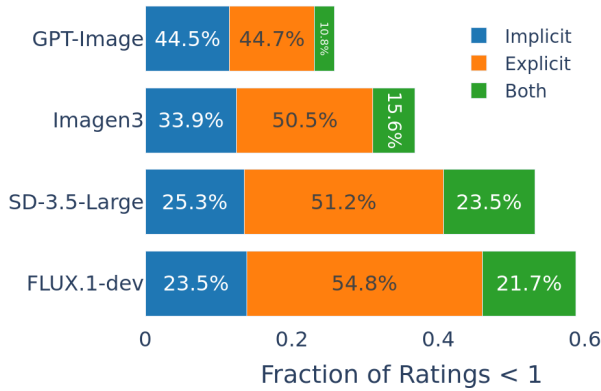


Figure 4: Distribution of image-prompt alignment errors (score < 1) by model, grouped by error type: implicit, explicit, or both. Bar lengths show fraction of total errors; % show each type’s share of model’s total errors.

higher total error rates where explicit errors form the largest share of their respective alignment failures. These results indicate that improvements are needed in both explicit and implicit cultural modeling.

In Canada, Poland, Germany, and Brazil, approximately two-thirds of comments mention explicit prompt mismatches, indicating that literal fidelity dominates their feedback. Conversely, annotator feedback from India, China, and South Africa is more evenly distributed, with roughly half of the remarks targeting explicit flaws and half targeting implicit cultural elements. At the opposite end of the spectrum, annotators from Japan and Iran predominantly highlight implicit cultural elements, such as absent rituals, attire, or local setting, with only about one-third of their comments citing explicit tokens. Chile follows the latter trend, albeit less strongly. Collectively, these observations indicate that T2I models increasingly fail to capture users’ implicit cultural expectations in regions like Asia and the Middle East, as contrasted with user feedback from the Americas and Europe.

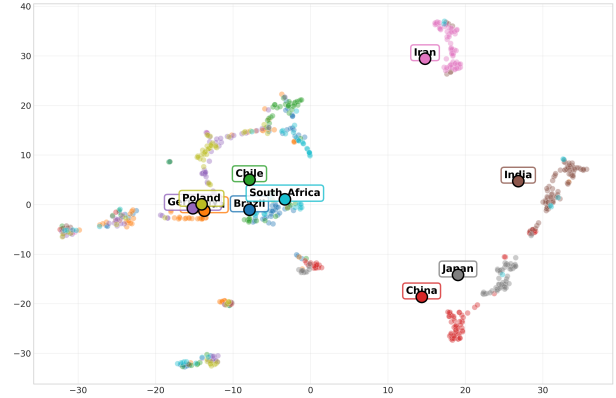


Figure 5: tSNE plot of Imagen3 images. Labeled markers show image embedding centroids per country.

Which words do models most frequently misinterpret?

Fig. 15 displays every word in the prompt that at least one rater labeled as erroneous, revealing two striking patterns. First, country demonyms (e.g., Iranian, Brazilian, Chinese, Japanese) are prominent. A closer examination of the rater comments reveals these words are typically highlighted as errors for two reasons: (i) a country-specific element is missing from the image, or (ii) the annotators are not able to relate to the depicted content. Second, terms such as *family*, *festival*, *ceremony*, *wedding*, *temple*, *meal*, *guests*, *tea*, *greeting*, *music*, *costumes*, and *flags* account for much of the remaining error frequency. These words represent broad cultural signifiers—rituals, social roles, and iconic objects—indicating that T2I models frequently misrepresent such elements.

What are the main causes of model failures across different countries?

To identify reasons behind model failures, we analyze free-form comments collected from annotators. For each country, we embed the comments using a sentence transformer⁵ and cluster them using HDBScan (Campello

⁵<https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

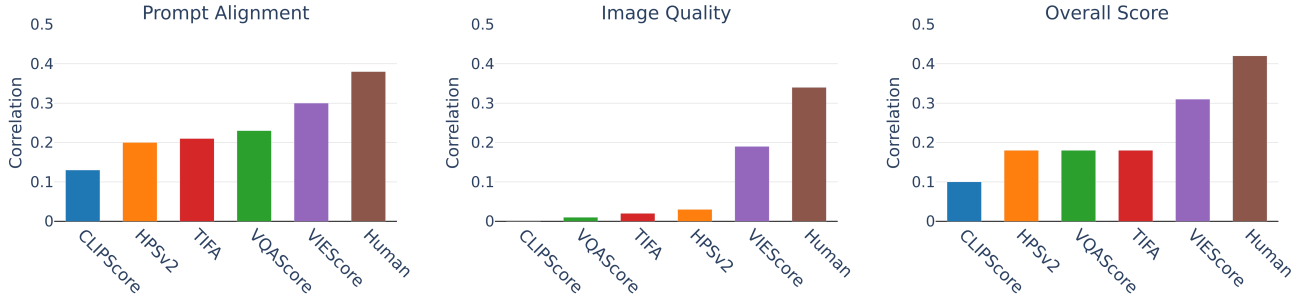


Figure 6: Spearman rank correlation of various T2I evaluation metrics with human ratings across three criteria: prompt alignment, image quality, and overall score. Human denotes the human-human Spearman rank correlation.

et al., 2013). We then prompt GPT-4o to summarize each cluster with a concise label and explanations. This approach reveals distinct failure patterns across regions. In Asia, models frequently misrepresent traditions and religious practices, often relying on stereotypes. In African contexts, outputs lacked cultural authenticity, defaulting to generic or Westernized portrayals. South American outputs suffered from poor regional specificity and inaccurate depictions of people’s appearances. Similarly, German outputs are consistently marked by stereotypical associations; Canadian content lacked appropriate demographic diversity and Indigenous representation. Further, we investigate the nature of the generated images by embedding them using the CLIP vision encoder.⁶ As shown in Fig. 5, images from Asian countries form distinct clusters, while those from other regions lack such clear grouping. This suggests model outputs fail to capture culturally distinctive visuals, demonstrating that failures are not uniform but potentially reflect specific training data blind spots and uneven geo-cultural representation.

6. Evaluating T2I Metrics on CULTURALFRAMES

Metrics analyzed. We analyze five representative metrics spanning different evaluation paradigms: CLIPScore (Hessel et al., 2021), TIFA (Hu et al., 2023), HPSv2 (Wu et al., 2023a), VQAScore (Lin et al., 2025), and VIEScore (Ku et al., 2024b). For TIFA, we use GPT-4o-mini as the question generation model and Qwen2.5-VL-32B-Instruct (Team, 2025) as the VQA module. GPT-4o is also used as the backbone VLM in VIEScore.

How do metrics perform against different rating criteria? We evaluate how well current T2I metrics correlate with human judgments across prompt alignment, image quality, and overall score (see Fig. 6). Among the evaluated

metrics, VIEScore achieves the highest correlation with human ratings across all criteria. For *prompt alignment*, VIEScore attains a Spearman correlation of 0.30. While this is below the human-human agreement of 0.38, it notably outperforms all other metrics. In contrast, TIFA, despite being explicitly designed to assess image-text faithfulness, exhibits a lower correlation, highlighting a gap between metric design and actual alignment with human perception. The performance gap is even more pronounced for *image quality*, where all metrics correlate poorly with human ratings. Nevertheless, VIEScore again performs best, followed by HPSv2. The relatively stronger performance of HPSv2 may be attributed to its alignment on image pairs, with human preference likely driven by image quality, potentially making it more sensitive to visual appeal. However, the overall weak correlations suggest that current metrics fail to capture the subjective nature of image quality as assessed by humans. For the *overall score*, VIEScore again demonstrates the highest alignment with human judgments, achieving a correlation of 0.31 compared to human-human agreement of 0.42. CLIPScore, in contrast, consistently underperforms, indicating limitations as a general-purpose evaluation metric, particularly for culturally sensitive image assessments.

Do explanations provided by VLM-based metrics capture the mistakes human raters highlight? To further analyze the effectiveness of the best-performing metric on our benchmark, VIEScore, we evaluate whether its generated explanations reflect the issues raised by human annotators. We adopt an LLM-as-a-judge setting, instructing it to assess the alignment between VIEScore’s reasoning and human concerns on a 1–5 Likert scale. The instructions are shown in Fig. 16. To calibrate the LLM’s judgments, we provided five in-context examples corresponding to varying quality levels. Additionally, we manually evaluate 100 judge-provided scores, sampled across countries and rating categories. We confirm that the LLM judge provides high-quality assessments. The results reveal that VIEScore’s explanations achieve an average rating of 2.19, indicating

⁶<https://huggingface.co/openai/clip-vit-large-patch14>

that while some overlap exists, the metric only partially captures the concerns raised by human raters. This also suggests a mismatch in the underlying rationale, emphasizing that current metrics, have substantial room for improvement in aligning with human judgment and reasoning processes. Some qualitative examples are provided in Fig. 17.

Can we improve metric performance through explicit instructions? Current T2I metrics are not explicitly guided to consider implicit and explicit prompt elements when evaluating image alignment. To test whether such guidance improves performance, we modify the instructions given to GPT-4o within VIEScore, replacing them with the more detailed annotation guidelines provided to human raters, including illustrative examples. We then re-evaluate images for image-prompt alignment using this instruction-tuned version of the VIEScore. This intervention yields a modest improvement in correlation with human ratings, with the Spearman correlation increasing from 0.30 to 0.32. To assess whether the reasoning behind the scores also improved, we again use the LLM-as-judge setup to evaluate 100 generated explanations. The resulting average score of 2.37, compared to 2.19 for the original VIEScore explanations, suggests that the modified metric captures human concerns slightly more effectively. Despite this improvement, the metric’s reasoning still falls considerably short of human rationale, indicating that explicit instructions alone are insufficient. These results underscore a persistent cultural and conceptual gap in model reasoning, even when provided with explicit guidance.

7. Conclusions

In this work, we introduce CULTURALFRAMES, a novel benchmark comprising 983 cultural prompts, 3,637 generated images, and 10,911 human annotations, spanning ten countries and five socio-cultural domains. CULTURALFRAMES assesses the ability of T2I models to generate images across diverse cultural contexts. We find that state-of-the-art T2I models not only fail to meet the more nuanced implicit expectations, but also the less challenging explicit expectations. In fact, models fail to meet cultural expectations 44% of the time on average across countries. Failures to meet explicit expectations averaged a surprisingly high 68% across models and countries, with implicit expectation failures also significant at 49%. Finally, we demonstrate that existing T2I evaluation metrics correlate poorly with human judgments of cultural alignment.

8. Limitations

Our study faces limitations due to our data collection methods and the scope of the CULTURALFRAMES. We approximated cultural groups as countries for annotator recruitment,

which may potentially oversimplify cultural identities and conflate culture with nationality due to practical constraints like information available in CulturalAtlas and annotator availability.

Our strategic choice to maximize diversity by recruiting multiple annotators per country, while enriching the evaluation with varied viewpoints, inherently presents a trade-off. A broader range of interpretations, stemming from a more diverse group, can naturally lead to lower inter-rater agreement scores when compared to evaluations conducted by a smaller, more homogenous annotator pool. It is this trade-off, coupled with the inherent subjectivity of the task, that provides context for our inter-annotator agreement results. This reflects the inherent subjectivity of evaluating cultural nuances and expectations.

A further limitation, driven by practical considerations of scale, is a generation of only a single image per model for each prompt. This single-instance evaluation makes it challenging for annotators to definitively identify stereotypical associations, as patterns of representation across multiple generations for the same prompt cannot be observed.

9. Ethical Considerations

Our CULTURALFRAMES benchmark comprises prompts and generated images, whose cultural alignment is rated by professional annotators via Prolific from the relevant countries. To ensure wide cultural representation, we recruited annotators from three distinct community groups within these countries, compensating them at \$10-15 per hour for all tasks performed, a rate established after pilot testing. This reflects our commitment to fair and inclusive data collection practices.

Despite the efforts, we acknowledge a key limitation: equating cultural groups with national borders within or across these national lines. This simplification may overlook the complex realities of minority and diaspora communities. We thus urge future research to explore finer-grained distinctions within cultural groups. While recognizing these constraints, we are hopeful that our work contributes to a deeper understanding of cultural nuances in visual generations and provides a foundation for such future investigations.

References

- Adilazuarda, M. F., Mukherjee, S., Lavania, P., Singh, S., Aji, A. F., O’Neill, J., Modi, A., and Choudhury, M. Towards measuring and modeling “culture” in LLMs: A survey, 2024. URL <https://arxiv.org/abs/2403.15412>.
- Akbulut, C., Robinson, K., Rauh, M., Albuquerque, I., Wiles, O., Weidinger, L., Rieser, V., Hasson, Y., Marchal,

- N., Gabriel, I., Isaac, W., and Hendricks, L. A. Century: A framework and dataset for evaluating historical contextualisation of sensitive images. In *International Conference on Learning Representations (ICLR)*, 2025. URL <https://openreview.net/forum?id=1KLBvrYz3V>.
- Alsudais, A. Analyzing how text-to-image models represent nationalities in everyday tasks, 2025. URL <https://arxiv.org/abs/2504.06313>.
- Basu, A., Babu, R. V., and Pruthi, D. Inspecting the geographical representativeness of images from text-to-image models, 2023. URL <https://arxiv.org/abs/2305.11080>.
- Bayramli, Z., Suleymanzade, A., An, N. M., Ahmad, H., Kim, E., Park, J., Thorne, J., and Oh, A. Diffusion models through a global lens: Are they culturally inclusive?, 2025. URL <https://arxiv.org/abs/2502.08914>.
- Bird, C., Ungless, E. L., and Kasirzadeh, A. Typology of risks of generative text-to-image models, 2023. URL <https://arxiv.org/abs/2307.05543>.
- Campello, R. J., Moulavi, D., and Sander, J. Density-based clustering based on hierarchical density estimates. In *Pacific-Asia conference on knowledge discovery and data mining*, pp. 160–172. Springer, 2013. URL <https://portal.findresearcher.sdu.dk/en/publications/density-based-clustering-based-on-hierarchical-density-estimates>.
- Cho, J., Zala, A., and Bansal, M. Dall-eval: Probing the reasoning skills and social biases of text-to-image generation models, 2023. URL <https://arxiv.org/abs/2202.04053>.
- Cho, J., Hu, Y., Garg, R., Anderson, P., Krishna, R., Baldridge, J., Bansal, M., Pont-Tuset, J., and Wang, S. Davidsonian scene graph: Improving reliability in fine-grained evaluation for text-to-image generation, 2024. URL <https://arxiv.org/abs/2310.18235>.
- Dehouche, N. and Dehouche, K. What’s in a text-to-image prompt? the potential of stable diffusion in visual arts education. *Heliyon*, 9(6):e16757, 2023. ISSN 2405-8440. doi: <https://doi.org/10.1016/j.heliyon.2023.e16757>. URL <https://www.sciencedirect.com/science/article/pii/S2405844023039646>.
- Esser, P., Kulal, S., Blattmann, A., Entezari, R., Müller, J., Saini, H., Levi, Y., Lorenz, D., Sauer, A., Boesel, F., Podell, D., Dockhorn, T., English, Z., Lacey, K., Goodwin, A., Marek, Y., and Rombach, R. Scaling rectified flow transformers for high-resolution image synthesis, 2024. URL <https://arxiv.org/abs/2403.03206>.
- Haerpfer, C., Inglehart, R., Moreno, A., Welzel, C., Kizilova, K., Lagos, M., Diez-Medrano, J., Norris, P., Ponarin, E., and Puranen, B. World Values Survey: Round seven - country-pooled datafile version 3.0. Madrid, Spain & Vienna, Austria: JD Systems Institute & WVSA Secretariat, 2022. URL <https://www.worldvaluessurvey.org/WVSDocumentationWV7.jsp>.
- Hall, M., Ross, C., Williams, A., Carion, N., Drozdal, M., and Soriano, A. R. Dig in: Evaluating disparities in image generations with indicators for geographic diversity, 2024. URL <https://arxiv.org/abs/2308.06198>.
- Hall, S. Encoding/decoding. In Hall, S., Hobson, D., Lowe, A., and Willis, P. (eds.), *Culture, Media, Language: Working Papers in Cultural Studies*, pp. 63–87. Hutchinson, London, 1980.
- Hartmann, J., Exner, Y., and Domdey, S. The power of generative marketing: Can generative ai create superhuman visual marketing content? *International Journal of Research in Marketing*, 42(1):13–31, 2025. ISSN 0167-8116. doi: <https://doi.org/10.1016/j.ijresmar.2024.09.002>. URL <https://www.sciencedirect.com/science/article/pii/S0167811624000843>.
- Herscovich, D., Frank, S., Lent, H., de Lhoneux, M., Abdou, M., Brandl, S., Bugliarello, E., Cabello Piqueras, L., Chalkidis, I., Cui, R., Fierro, C., Margatina, K., Rust, P., and Søgaard, A. Challenges and strategies in cross-cultural NLP. In Muresan, S., Nakov, P., and Villavicencio, A. (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 6997–7013, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.482. URL <https://aclanthology.org/2022.acl-long.482>.
- Hessel, J., Holtzman, A., Forbes, M., Le Bras, R., and Choi, Y. CLIPScore: A reference-free evaluation metric for image captioning. In Moens, M.-F., Huang, X., Specia, L., and Yih, S. W.-t. (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 7514–7528, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.595. URL <https://aclanthology.org/2021.emnlp-main.595/>.
- Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S. Gans trained by a two time-scale update rule converge to a local nash equilibrium, 2018. URL <https://arxiv.org/abs/1706.08500>.

- Hofstede, G., Hofstede, G. J., and Minkov, M. *Cultures and organizations: software of the mind: intercultural cooperation and its importance for survival*. McGraw-Hill, New York; London, 3rd edition, 2010. URL <https://www.mhprofessional.com/cultures-and-organizations-software-of-the-mind-third-edition-9780071664189-usa>.
- Hu, Y., Liu, B., Kasai, J., Wang, Y., Ostendorf, M., Krishna, R., and Smith, N. A. Tifa: Accurate and interpretable text-to-image faithfulness evaluation with question answering, 2023. URL <https://arxiv.org/abs/2303.11897>.
- Huang, K., Sun, K., Xie, E., Li, Z., and Liu, X. T2i-compbench: A comprehensive benchmark for open-world compositional text-to-image generation. *Advances in Neural Information Processing Systems*, 36:78723–78747, 2023.
- Huang, K., Duan, C., Sun, K., Xie, E., Li, Z., and Liu, X. T2i-compbench++: An enhanced and comprehensive benchmark for compositional text-to-image generation, 2025. URL <https://arxiv.org/abs/2307.06350>.
- Imagen-Team-Google. Imagen 3, 2024. URL <https://arxiv.org/abs/2408.07009>.
- Jeong, S., Choi, I., Yun, Y., and Kim, J. Culture-trip: Culturally-aware text-to-image generation with iterative prompt refinement, 2025. URL <https://arxiv.org/abs/2502.16902>.
- Jha, A., Prabhakaran, V., Denton, R., Laszlo, S., Dave, S., Qadri, R., Reddy, C., and Dev, S. ViSAGE: A global-scale analysis of visual stereotypes in text-to-image generation. In Ku, L.-W., Martins, A., and Srikumar, V. (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 12333–12347, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.667. URL <https://aclanthology.org/2024.acl-long.667/>.
- Jiang, D., Ku, M., Li, T., Ni, Y., Sun, S., Fan, R., and Chen, W. Genai arena: An open evaluation platform for generative models, 2024. URL <https://arxiv.org/abs/2406.04485>.
- Kannen, N., Ahmad, A., Andreetto, M., Prabhakaran, V., Prabhu, U., Dieng, A. B., Bhattacharyya, P., and Dave, S. Beyond aesthetics: Cultural competence in text-to-image models, 2025. URL <https://arxiv.org/abs/2407.06863>.
- Khanuja, S., Ramamoorthy, S., Song, Y., and Neubig, G. An image speaks a thousand words, but can everyone listen? on image transcreation for cultural relevance. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 10258–10279, 2024.
- Kirstain, Y., Polyak, A., Singer, U., Matiana, S., Penna, J., and Levy, O. Pick-a-pic: An open dataset of user preferences for text-to-image generation, 2023. URL <https://arxiv.org/abs/2305.01569>.
- Krippendorff, K. *Content Analysis: An Introduction to Its Methodology*. SAGE Publications, 2013. ISBN 9781412983150. URL https://books.google.ch/books?id=s_yqFXnGgjQC.
- Ku, M., Jiang, D., Wei, C., Yue, X., and Chen, W. VIEScore: Towards explainable metrics for conditional image synthesis evaluation. In Ku, L.-W., Martins, A., and Srikumar, V. (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 12268–12290, Bangkok, Thailand, August 2024a. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.663. URL <https://aclanthology.org/2024.acl-long.663/>.
- Ku, M., Jiang, D., Wei, C., Yue, X., and Chen, W. VIEScore: Towards explainable metrics for conditional image synthesis evaluation. In Ku, L.-W., Martins, A., and Srikumar, V. (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 12268–12290, Bangkok, Thailand, August 2024b. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.663. URL <https://aclanthology.org/2024.acl-long.663/>.
- Ku, M., Li, T., Zhang, K., Lu, Y., Fu, X., Zhuang, W., and Chen, W. Imagenhub: Standardizing the evaluation of conditional image generation models, 2024c. URL <https://arxiv.org/abs/2310.01596>.
- Labs, B. F. Flux. <https://github.com/black-forest-labs/flux>, 2024.
- Lee, T., Yasunaga, M., Meng, C., Mai, Y., Park, J. S., Gupta, A., Zhang, Y., Narayanan, D., Teufel, H. B., Bellagente, M., Kang, M., Park, T., Leskovec, J., Zhu, J.-Y., Fei-Fei, L., Wu, J., Ermon, S., and Liang, P. Holistic evaluation of text-to-image models, 2023. URL <https://arxiv.org/abs/2311.04287>.
- Lin, Z., Pathak, D., Li, B., Li, J., Xia, X., Neubig, G., Zhang, P., and Ramanan, D. Evaluating text-to-visual generation with image-to-text generation. In Leonardis, A., Ricci, E., Roth, S., Russakovsky, O., Sattler, T., and

- Varol, G. (eds.), *Computer Vision – ECCV 2024*, pp. 366–384, Cham, 2025. Springer Nature Switzerland. ISBN 978-3-031-72673-6.
- Liu, Z., Schaldenbrand, P., Okogwu, B.-C., Peng, W., Yun, Y., Hundt, A., Kim, J., and Oh, J. Scoft: Self-contrastive fine-tuning for equitable image generation, 2024. URL <https://arxiv.org/abs/2401.08053>.
- Loukili, S., Elaachak, L., and Fennan, A. Finetuning stable diffusion models for email marketing text-to-image generation. In Ben Ahmed, M., Abdelhakim, B. A., Karas, I. R., and Ben Ahmed, K. (eds.), *Innovations in Smart Cities Applications Volume 8*, pp. 524–535, Cham, 2025. Springer Nature Switzerland. URL https://doi.org/10.1007/978-3-031-88653-9_51.
- Lu, Y., Yang, X., Li, X., Wang, X. E., and Wang, W. Y. LLM-Score: Unveiling the power of large language models in text-to-image synthesis evaluation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=OJ0c6um1An>.
- Maharana, A., Hannan, D., and Bansal, M. Storydall-e: Adapting pretrained text-to-image transformers for story continuation. In Avidan, S., Brostow, G., Cissé, M., Farinella, G. M., and Hassner, T. (eds.), *Computer Vision – ECCV 2022*, pp. 70–87, Cham, 2022. Springer Nature Switzerland. URL https://doi.org/10.1007/978-3-031-19836-6_5.
- McLuhan, M. *Understanding Media: The Extensions of Man*. Signet Books, New York, 1966.
- Mosaica. The cultural atlas. <https://culturalatlases.sbs.com.au/>, 2024.
- Naik, R. and Nushi, B. Social biases through the text-to-image generation lens, 2023. URL <https://arxiv.org/abs/2304.06034>.
- OpenAI. Gpt-4o system card, 2024. URL <https://arxiv.org/abs/2410.21276>.
- OpenAI. Introducing 4o image generation. <https://openai.com/index/introducing-4o-image-generation/>, 2025.
- Prabhakaran, V., Qadri, R., and Hutchinson, B. Cultural incongruencies in artificial intelligence, 2022. URL <https://arxiv.org/abs/2211.13069>.
- Qadri, R., Diaz, M., Wang, D., and Madaio, M. The case for “thick evaluations” of cultural representation in AI, 2025. URL <https://arxiv.org/abs/2503.19075>.
- Rastogi, C., Teh, T. H., Mishra, P., Patel, R., Ashwood, Z., Davani, A. M., Diaz, M., Paganini, M., Parrish, A., Wang, D., Prabhakaran, V., Aroyo, L., and Rieser, V. Insights on disagreement patterns in multimodal safety perception across diverse rater groups, 2024. URL <https://arxiv.org/abs/2410.17032>.
- Ruiz, N., Li, Y., Jampani, V., Pritch, Y., Rubinstein, M., and Aberman, K. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation, 2023. URL <https://arxiv.org/abs/2208.12242>.
- Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E., Ghasemipour, S. K. S., Ayan, B. K., Mahdavi, S. S., Lopes, R. G., Salimans, T., Ho, J., Fleet, D. J., and Norouzi, M. Photorealistic text-to-image diffusion models with deep language understanding, 2022. URL <https://arxiv.org/abs/2205.11487>.
- Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., and Chen, X. Improved techniques for training gans, 2016. URL <https://arxiv.org/abs/1606.03498>.
- Struppek, L., Hintersdorf, D., Friedrich, F., Br, M., Schramowski, P., and Kersting, K. Exploiting cultural biases via homoglyphs in text-to-image synthesis. *Journal of Artificial Intelligence Research*, 78: 1017–1068, December 2023. ISSN 1076-9757. doi: 10.1613/jair.1.15388. URL <http://dx.doi.org/10.1613/jair.1.15388>.
- Team, G. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context, 2024. URL <https://arxiv.org/abs/2403.05530>.
- Team, Q. Qwen2.5-vl, January 2025. URL <https://qwenlm.github.io/blog/qwen2.5-vl/>.
- Ventura, M., Ben-David, E., Korhonen, A., and Reichart, R. Navigating cultural chasms: Exploring and unlocking the cultural pov of text-to-image models, 2024. URL <https://arxiv.org/abs/2310.01929>.
- Wan, Y., Subramonian, A., Ovalle, A., Lin, Z., Suvama, A., Chance, C., Bansal, H., Pattichis, R., and Chang, K.-W. Survey of bias in text-to-image generation: Definition, evaluation, and mitigation, 2024. URL <https://arxiv.org/abs/2404.01030>.
- Wu, X., Hao, Y., Sun, K., Chen, Y., Zhu, F., Zhao, R., and Li, H. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis, 2023a. URL <https://arxiv.org/abs/2306.09341>.

- Wu, X., Sun, K., Zhu, F., Zhao, R., and Li, H. Human preference score: Better aligning text-to-image models with human preference. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2096–2105, 2023b.
- Xu, J., Liu, X., Wu, Y., Tong, Y., Li, Q., Ding, M., Tang, J., and Dong, Y. Imagereward: Learning and evaluating human preferences for text-to-image generation, 2023. URL <https://arxiv.org/abs/2304.05977>.
- Yarom, M., Bitton, Y., Changpinyo, S., Aharoni, R., Herzig, J., Lang, O., Ofek, E., and Szpektor, I. What you see is what you read? improving text-image alignment evaluation. In *NeurIPS*, 2023.
- Yerukola, A., Gabriel, S., Peng, N., and Sap, M. Mind the gesture: Evaluating ai sensitivity to culturally offensive non-verbal gestures, 2025. URL <https://arxiv.org/abs/2502.17710>.
- Yu, J., Xu, Y., Koh, J. Y., Luong, T., Baid, G., Wang, Z., Vasudevan, V., Ku, A., Yang, Y., Ayan, B. K., Hutchinson, B., Han, W., Parekh, Z., Li, X., Zhang, H., Baldrige, J., and Wu, Y. Scaling autoregressive models for content-rich text-to-image generation, 2022. URL <https://arxiv.org/abs/2206.10789>.

A. Appendix

A.1. CULTURALFRAMES

This section outlines the full pipeline used to create the CULTURALFRAMES. We describe how culturally grounded prompts were generated, filtered, and verified by human annotators across multiple countries. We also detail how these prompts were used to generate images from various text-to-image models, along with the settings and parameters used for generation.

A.1.1. PROMPT GENERATION

We begin with the Cultural Atlas ([Mosaica, 2024](#)), a curated knowledge base of cross-cultural attitudes, practices, norms, behaviors, and communication styles, designed to inform and educate the public about Australia’s migrant populations. The Atlas provides detailed textual descriptions across categories such as family structures, greeting customs, cultural etiquette, religious beliefs, and more. We use the Cultural Atlas as a source of culturally grounded information to guide prompt generation. However, not all categories in the Atlas are suitable for visual depiction. We selected five categories—*dates-of-significance*, *etiquette*, *family*, *religion*, and *greetings*—based on two main criteria: (1) the content describes values or practices that can be meaningfully represented in images, and (2) the category is consistently available across a broad set of countries to support cross-cultural comparison.

We parsed the textual content from each selected category and segmented it into paragraphs using newline characters. Each paragraph served as an input “excerpt” to an LLM for prompt generation. Given a country and an excerpt, we prompted GPT-4o (gpt-4o-2024-08-06) ([OpenAI, 2024](#)) to generate two short prompts (each under 15 words) that: (i) were grounded in the excerpt’s content, (ii) described a culturally relevant and visually observable scenario, and (iii) included sufficient country-specific context, either explicitly or implicitly. The prompts were designed to reflect underlying cultural values through everyday, observable situations, such as a wedding ceremony or a workplace interaction. To guide this process, we crafted category-specific instructions that encouraged the model to generate meaningful and culturally grounded prompts.

We began by generating a small number of prompts per category, which were evaluated by human annotators to assess whether the scenarios were both visually depictable and culturally appropriate (see Section A.1.2 for details). Prompts that passed these quality checks were reused as few-shot in-context examples to guide further prompt generation. This iterative process enabled us to scale prompt creation while maintaining cultural fidelity and diversity. Instructions provided to GPT-4o ([OpenAI, 2024](#)) used across different

Country	Unique Annotators	Avg Age	% Male	% Female	% Other
Brazil	35	36.1	69.0	31.0	0.0
Canada	34	37.9	47.9	52.1	0.0
Chile	35	31.1	77.7	22.3	0.0
China	40	33.0	32.3	67.7	0.0
Germany	51	35.1	68.5	31.5	0.0
India	32	31.7	46.6	53.4	0.0
Iran	28	32.0	47.0	53.0	0.0
Japan	25	44.2	56.1	40.6	3.2
Poland	27	32.0	62.0	38.0	0.0
South Africa	83	32.9	35.1	64.9	0.0

Table 2: Summary of participant demographics by country.

categories are provided below.

A.1.2. PROMPT FILTERING

For every country, we ask 3 culturally knowledgeable annotators if the prompt represents a scenario observable in their culture and aligns with their values. Only those prompts that 2 or more annotators choose make it into CULTURALFRAMES. In Fig. 7, we present the prompt filtering interface where annotators choose “Yes/No” for a given prompt depending on whether the prompt reflects an observable scenario in their culture that aligns with their cultural values.

A.1.3. PROMPT DISTRIBUTION ACROSS CATEGORIES

Fig. 8 shows the distribution of prompts across five cultural categories used in constructing CULTURALFRAMES: *dates-of-significance*, *etiquette*, *family*, *religion*, and *greetings*. Across countries, *dates-of-significance* consistently accounts for the largest share of prompts, followed by *etiquette*. This distribution reflects the relative amount of information available for each category in the Cultural Atlas. The remaining three categories—*family*, *religion*, and *greetings*—have relatively balanced proportions. We aimed to maintain a similar category distribution across countries to support fair cross-cultural comparisons. Notably, South Africa lacks sufficient information in the *family* category, so it is excluded from that category in the figure.

A.1.4. IMAGE GENERATION

We generate images at a resolution of 1024×1024 across all models to ensure consistency. For GPT-Image, we set the image quality to high. For Imagegen3, we use VertexAI to make API calls and enable the default enhance_prompt setting, which expands the prompt prior to image generation. For FLUX.1-dev, we set the guidance scale to 3.5, max_sequence_length to 512, and use 50 inference steps. In the case of SD-3.5-Large, we use a guidance scale of 4.5

and 40 inference steps.

A.2. Image Rating

We develop a custom interface for collecting image ratings. Fig. 10 and Fig. 11 show the detailed instructions we provide to the annotators for rating images. Fig. 12 shows the interface where annotators rate images.

A.2.1. ANNOTATOR DEMOGRAPHICS

Tab. 2 provides details on the annotators who participated in our studies.

Prompt Instructions (Religion)**Purpose:**

We want to test whether text-to-image models can accurately capture how religion is practiced in a particular country along with its norms, practices, rituals, traditions, and values. You will be given:

1. A country name
2. A short excerpt on religious norms: an implicit description of how religion is practiced or influences everyday life, or some information that is related to religious practices.

Your Task:

Use these inputs to produce two short prompts (each under 15 words) that is rooted in the provided excerpt and explore diverse scenarios, to evaluate the image-generation model's understanding of the religion of the country. Each prompt should:

- Be clearly rooted in the excerpt's details and context (e.g., setting, participants, timing). You must not deviate from the provided excerpt
- Create prompts that describe specific daily interactions, rituals, or scenarios that reflect the cultural values and social norms related to religion and mentioned in the excerpt. These should be concrete, observable situations that commonly occur in this culture/country.
- Be diverse, realistic scenario, and under 15 words
- Be visually depictable - that is, it must be possible to generate a meaningful and culturally relevant image based on the prompt.

Important: Make sure the country can be inferred from the prompt. It should be either stated explicitly like mentioning a region or name of the country or there must be enough country specific elements in the prompt to infer the country.

Note: If the information provided cannot be used to create a practical observable scenario that can be depicted as an image, return "N/A".

Return the prompts in this JSON format:

```
{
  "prompt_1": "...",
  "prompt_2": "..."
}
```

Here are the inputs:

- Country: {country}
- Excerpt: {excerpt}

Previously Generated Prompts (to avoid duplication):

```
{already_generated_prompts}
```

Accepted Examples:

```
{incontext_examples_positive}
```

Rejected Examples:

```
{incontext_examples_negative}
```

Generate **exactly two** new prompts that satisfy all of the criteria above, follow the style/patterns of the accepted examples, avoid the issues shown in the rejected ones, and explore diverse scenarios different from the ones already generated. Output **only** the JSON object specified.

Prompt Instructions (Etiquette)**Purpose:**

We want to test whether text-to-image models can accurately capture how etiquette is practiced in a particular country, including norms, manners, and social conduct related to visiting, gifting, eating, and other social situations. You will be given:

1. A country name
2. A short excerpt on etiquette norms: an implicit description of how people in this country engage with each other in different social situations, or some information related to etiquette.

Your Task:

Use these inputs to produce two short prompts (each under 15 words) that is rooted in the provided excerpt and explore diverse scenarios, to evaluate the image-generation model's understanding of etiquette. Each prompt should:

- Be clearly rooted in the excerpt's details and context (e.g., setting, participants, timing). You must not deviate from the provided excerpt
- Represent a social scenario or interaction where the etiquette norm or value mentioned in the excerpt can be observed. It must be a realistic, observable scenario that commonly occurs in this culture/country.
- Do not explicitly name the etiquette rule. Be implicit in conveying the details. The goal is to create situations where the etiquette rule can be observed and inferred by the model.
- Be diverse, realistic scenario, and under 15 words
- Be visually depictable - that is, it must be possible to generate a meaningful and culturally relevant image based on the prompt.
- Avoid using phrases like "arriving late", "arriving on time" and other such phrases that cannot be visualized in the image.

Important: Make sure the country can be inferred from the prompt. It should be either stated explicitly like mentioning a region or name of the country or there must be enough country specific elements in the prompt to infer the country.

Note: If the information provided cannot be used to create a practical observable scenario that can be depicted as an image, return "N/A".

Return the prompts in this JSON format:

```
{
  "prompt_1": "...",
  "prompt_2": "..."
}
```

Here are the inputs:

- Country: {country}
- Excerpt: {excerpt}

Previously Generated Prompts (to avoid duplication):

```
{already_generated_prompts}
```

Accepted Examples:

```
{incontext_examples_positive}
```

Rejected Examples:

```
{incontext_examples_negative}
```

Generate **exactly two** new prompts that satisfy all of the criteria above, follow the style/patterns of the accepted examples, avoid the issues shown in the rejected ones, and explore diverse scenarios different from the ones already generated. Output **only** the JSON object specified.

Prompt Instructions (Family)

Purpose:

We want to test whether text-to-image models can accurately depict how family values, structures, and dynamics operate in a particular country. You will be given:

1. A country name
2. A short excerpt on family norms: an implicit description of how family life, roles, or relationships function in this culture.

Your Task:

Use these inputs to produce two short prompts (each under 12 words) that are clearly rooted in the provided excerpt and explore diverse scenarios, to evaluate a model’s understanding of these family practices. Each prompt should:

- Be firmly based on the excerpt’s context. You must not deviate from the provided excerpt
- Portray family related interactions that happen in the culture/country conditioned on the values, norms provided in the excerpt
- Avoid explicitly naming the core family norm or value, but include enough detail for the model to infer it
- Depict diverse, realistic scenarios that convey familial interactions, each under 12 words
- Be visually depictable - that is, it must be possible to generate a meaningful and culturally relevant image based on the prompt.

Important: Make sure the country can be inferred from the prompt. It should be either stated explicitly like mentioning a region or name of the country or there must be enough country specific elements in the prompt to infer the country.

Note: If the information provided cannot be used to create a practical observable scenario that can be depicted as an image, return "N/A".

Return the prompts in this JSON format:

```
{
  "prompt_1": "...",
  "prompt_2": "..."
}
```

Here are the inputs:

- Country: {country}
- Excerpt: {excerpt}

Previously Generated Prompts (to avoid duplication):

{already_generated_prompts}

Accepted Examples:

{incontext_examples_positive}

Rejected Examples:

{incontext_examples_negative}

Generate **exactly two** new prompts that satisfy all of the criteria above, follow the style/patterns of the accepted examples, avoid the issues shown in the rejected ones, and explore diverse scenarios different from the ones already generated. Output **only** the JSON object specified.

Prompt Instructions (Dates-of-significance)**Purpose:**

We want to test whether text-to-image models can accurately depict how a country observes its significant dates—festivals, holidays, or other notable events. You will be given:

1. A country name
2. A short excerpt on a date of significance: an implicit description of festivities, traditions, or commemorative practices related to this important day.

Your Task:

Use these inputs to produce two short prompts (under 12 words) that are clearly rooted in the provided excerpt and explore diverse scenarios, to evaluate a model's understanding of these celebrations. Each prompt should:

- Be firmly based on the excerpt's context. You must not deviate from the provided excerpt
- Represent daily interactions, rituals, or scenarios that are related to this date of significance. It must be a realistic, observable scenario that commonly occurs in this culture/country.
- Convey the date of significance through rituals, traditions, or celebrations that are specific to this date.
- Depict diverse, realistic scenarios that convey how people observe this date, each under 12 words.
- Be visually depictable - that is, it must be possible to generate a meaningful and culturally relevant image based on the prompt.

Important: Make sure the country can be inferred from the prompt. It should be either stated explicitly like mentioning a region or name of the country or there must be enough country specific elements in the prompt to infer the country.

Note: If the information provided cannot be used to create a practical observable scenario that can be depicted as an image, return "N/A".

Return the prompts in this JSON format:

```
{
  "prompt_1": "...",
  "prompt_2": "..."
}
```

Here are the inputs:

- Country: {country}
- Excerpt: {excerpt}

Previously Generated Prompts (to avoid duplication):

```
{already_generated_prompts}
```

Accepted Examples:

```
{incontext_examples_positive}
```

Rejected Examples:

```
{incontext_examples_negative}
```

Generate **exactly two** new prompts that satisfy all of the criteria above, follow the style/patterns of the accepted examples, avoid the issues shown in the rejected ones, and explore diverse scenarios different from the ones already generated. Output **only** the JSON object specified.

Prompt Validation

Prompt 1 of 10

Prompt:

American family dining, engaging in lively conversation while eating dinner

Does the prompt describe an observable scenario in your culture that aligns with your cultural values, norms, and practices and can be depicted as an image?

Yes

No

Continue

Figure 7: Prompt filtering interface where annotators choose “Yes/No” for a given prompt depending on whether the prompt reflects an observable scenario in their culture that aligns with their cultural values.



Figure 8: Distribution of prompts from different categories across countries.



Poppies worn on lapels during Remembrance Day ceremony

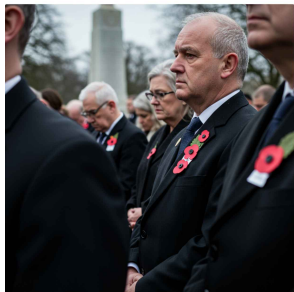
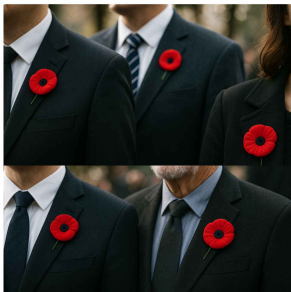


Imagen3



GPT-Image



SD-3.5 Large



FLUX.1.dev



Traditional Japanese tea ceremony in autumn garden

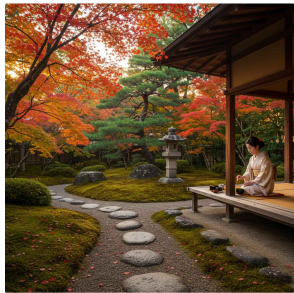
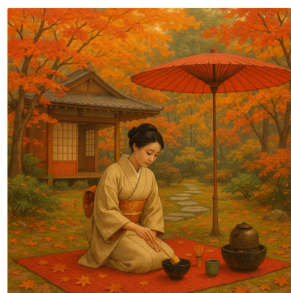


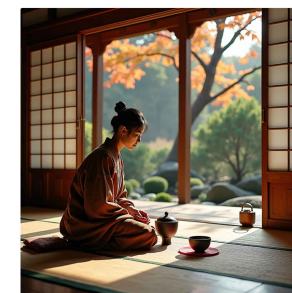
Imagen3



GPT-Image



SD-3.5 Large



FLUX.1.dev



Elderly Iranian man praying alone in a quiet mosque

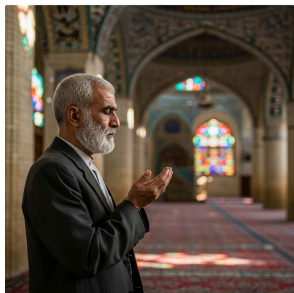
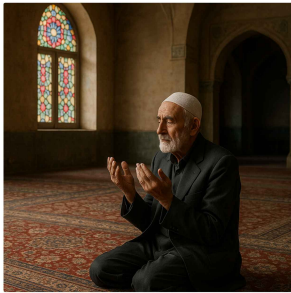
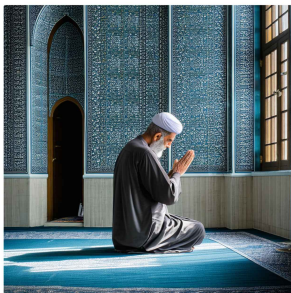


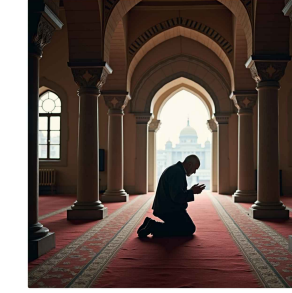
Imagen3



GPT-Image



SD-3.5 Large



FLUX.1.dev

Figure 9: Prompt-image examples from CULTURALFRAMES across different countries generated by the models.

Rating Criteria

You will rate each image on the following criteria:

1. Image-Prompt Alignment

Definition: You will evaluate how well the generated image matches the given prompt. You will assign a score of 0, 0.5, or 1 based on how faithful the generated image is with respect to the given text prompt.

What to look for: While evaluating the alignment, you should check for the faithfulness of the image with respect to both explicit and implicit elements in the prompt. See below for further details on explicit and implicit elements:

1. Explicit elements: These are elements clearly stated as words in the prompt, such as objects, actions, people, relationships, or settings. A good image must include all of these explicitly mentioned elements and represent them accurately.

Example of Explicit Elements



Prompt: "People offering flowers to Saraswati statue"

Here are the explicit elements in this prompt and how you can think about them:

- **People** - Are there any people in the image?
- **Offering** - Are the people offering something?
- **Flowers** - Are there any flowers in the image people are offering?
- **Saraswati statue** - Is there a Saraswati statue in the image?

For the image to align with the prompt, it must include all of these explicitly mentioned elements.

2. Implicit elements: These are elements of the prompt **that are not directly mentioned as words in the prompt** but are expected to be present in the image based on the cultural context. These may include appropriate attire or food for the setting, gestures or expressions that suit the context, interactions between people, or additional details that contribute to the authenticity of the scene. A strong image will reflect these expectations in a way that feels appropriate to someone familiar with the described scenario.

Example of Implicit Elements



Prompt: "People offering flowers to Saraswati statue"

Here are some implicit elements to look for and how you can think about them:

- **Setting** - Does the environment feel appropriate for a religious offering like a temple or home altar?
- **Attire** - Are the people dressed in a culturally appropriate way for the occasion?
- **Statue details** - Is the Saraswati statue depicted correctly with her common features, like a veena, white clothing, or a swan nearby?

These elements aren't directly mentioned in the prompt but are expected based on cultural context. You may notice others from your own experience. For the image to match the prompt, it should include and accurately show these details.

Meaning of the scores:

- **Score 1 (Good Match):** The image closely and accurately reflects the prompt. All key **explicit** elements and relevant **implicit** elements are present and appropriately depicted. The scene as a whole fits the situation described.
- **Score 0.5 (Partial Match):** The image shows a **reasonable attempt to follow the prompt**. Some explicit or implicit elements (or both) may be wrong or missing. However, the missing or incorrect elements are **not central to the scene** or **do not drastically change the intended meaning** of the prompt. Use this score when the image gets most of the prompt right, but lacks completeness.
- **Score 0 (Poor Match):** The image fails to meaningfully reflect the prompt. **Key explicit and/or implicit elements are missing or clearly incorrect**. Use this score when the image substantially misrepresents the described situation.

Note: You can access this information anytime during rating by clicking the ⓘ icon next to the criteria.

If you gave a score less than 1, please select what was missing or incorrect and provide a brief comment explaining what was wrong:

- Explicit elements (mentioned in the prompt) are missing or incorrect.
 - You **must** select the words in the prompt that are missing or incorrect and explain what was wrong.
- Implicit elements (not mentioned, but needed for cultural authenticity) are wrong or missing.
 - You can select the words that you think are missing or have incorrect implicit details, and you **must** explain what these missing or incorrect details are.

2. Stereotype

You will evaluate whether the image shows any stereotypes about your culture — that is, a narrow, over-generalized belief, or inaccurate view of it.

Answer with: **Yes** or **No**

If you have selected **Yes**, you must describe what aspect of the image is stereotypical.

3. Image Quality

Definition: You will evaluate the visual quality of the generated image. You will assign a score of 0, 0.5, or 1 based on whether the image looks natural, convincing, and contains any distortions or artifacts.

Meaning of the scores:

- **Score 1 (High Quality):** The image looks visually convincing and realistic. There are no visible distortions, artifacts, or unnatural elements. Objects, people and the scene are clear and harmoniously integrated.
- **Score 0.5 (Moderate Quality):** The image includes minor artifacts, distortions, or inconsistencies or, gives off an unnatural impression. However, most of the objects, people and the scene are still recognizable.
- **Score 0 (Poor Quality):** The image contains serious distortions, visual artifacts, or gives an unnatural impression or unusual sense that make objects or the scene hard to recognize or understand.

Note: You can access this information anytime during rating by clicking the ⓘ icon next to the criteria.

Artifacts and Unnatural Impression, respectively, are:

- **Artifacts:** Distortion, watermarks, scratches, blurred faces, unusual body parts (e.g., extra fingers, misshapen limbs), subjects not harmonized with the background
- **Unnatural Impression:** Wrong sense of distance (subject too big or too small compared to others), wrong shadows, incorrect lighting, unnatural colors, perspective issues

Examples (Click on the images to zoom in):



Score: 1

Clear image with natural proportions, good lighting, and no visible artifacts or distortions.



Score: 0.5

Minor distortions in facial features and unnaturally long hands, but overall scene is still recognizable.



Score: 0

Severe artifacts in hands with pig and hands morphed together making objects in the image difficult to recognise.

4. Overall Score

Definition: On a scale of 1 (very bad) to 5 (very good), how well do you think the image reflects the prompt?


Figure 11: Instructions given to annotators for stereotype, image quality, and overall score criteria.

Image 1 of 5

Prompt:

Chilean wedding, guests greeting bride's parents with warm handshakes

Click image to zoom in



Rate this Image

Image-Prompt Alignment ⓘ

Rate how well the image matches the given prompt.

0

0.5

1

Please select what was missing or incorrect:

☐ Explicit elements (mentioned in the prompt) are missing or incorrect

☐ Implicit elements (not mentioned, but needed for cultural authenticity) are wrong or missing

Select words from the prompt that weren't accurately depicted or missing:

Chileanweddingguestsgreetingbrideparentswarmhandshakes

Please explain what aspects were missing or incorrectly depicted:

Describe clearly what was missing or incorrect. This explanation is required for both explicit and implicit elements.

Stereotype

Does the image show any stereotypes about your culture — that is, a narrow, simplified, or inaccurate view of it?

No

Yes

Image Quality ⓘ

Rate the overall quality of the image, focusing on how clear, realistic, and natural it looks.

0

0.5

1

Overall Score

Rate how well do you think the image reflects the prompt

1

2

3

4

5

Figure 12: Rating collection interface shown to the annotators. When annotators select a score of less than 1, they need to give detailed feedback regarding explicit and implicit expectations, along with selecting the problematic words

23

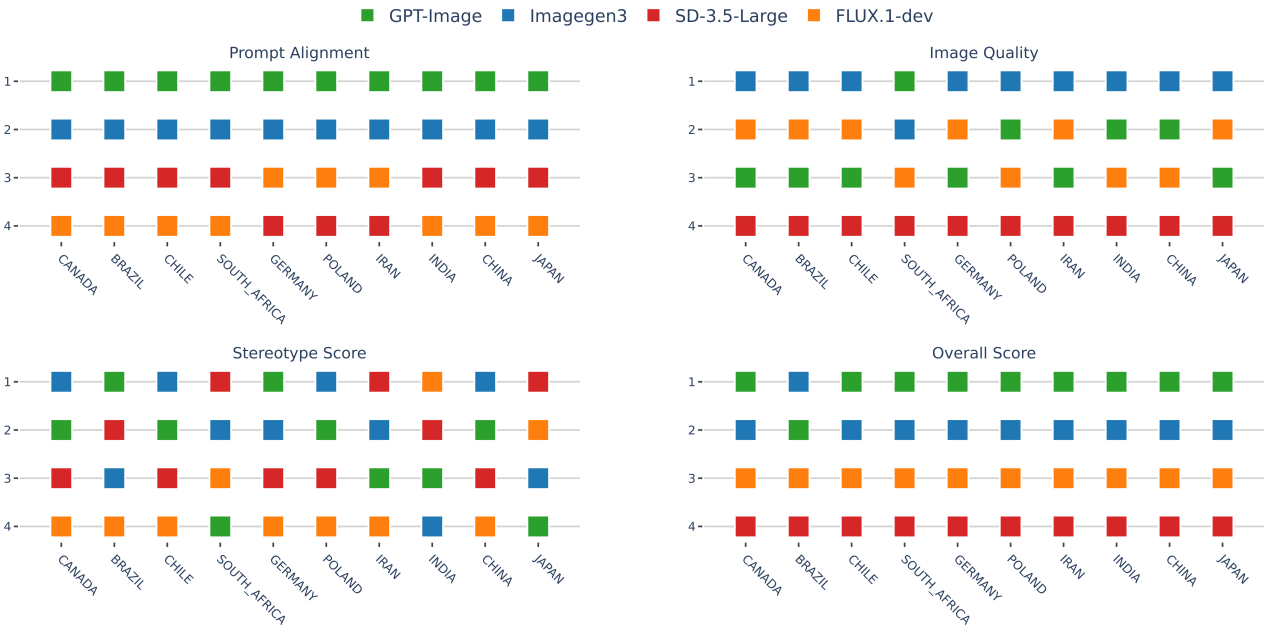


Figure 13: Model ranking across countries for different criteria (1 is the highest rank). Countries are grouped by geographical proximity.

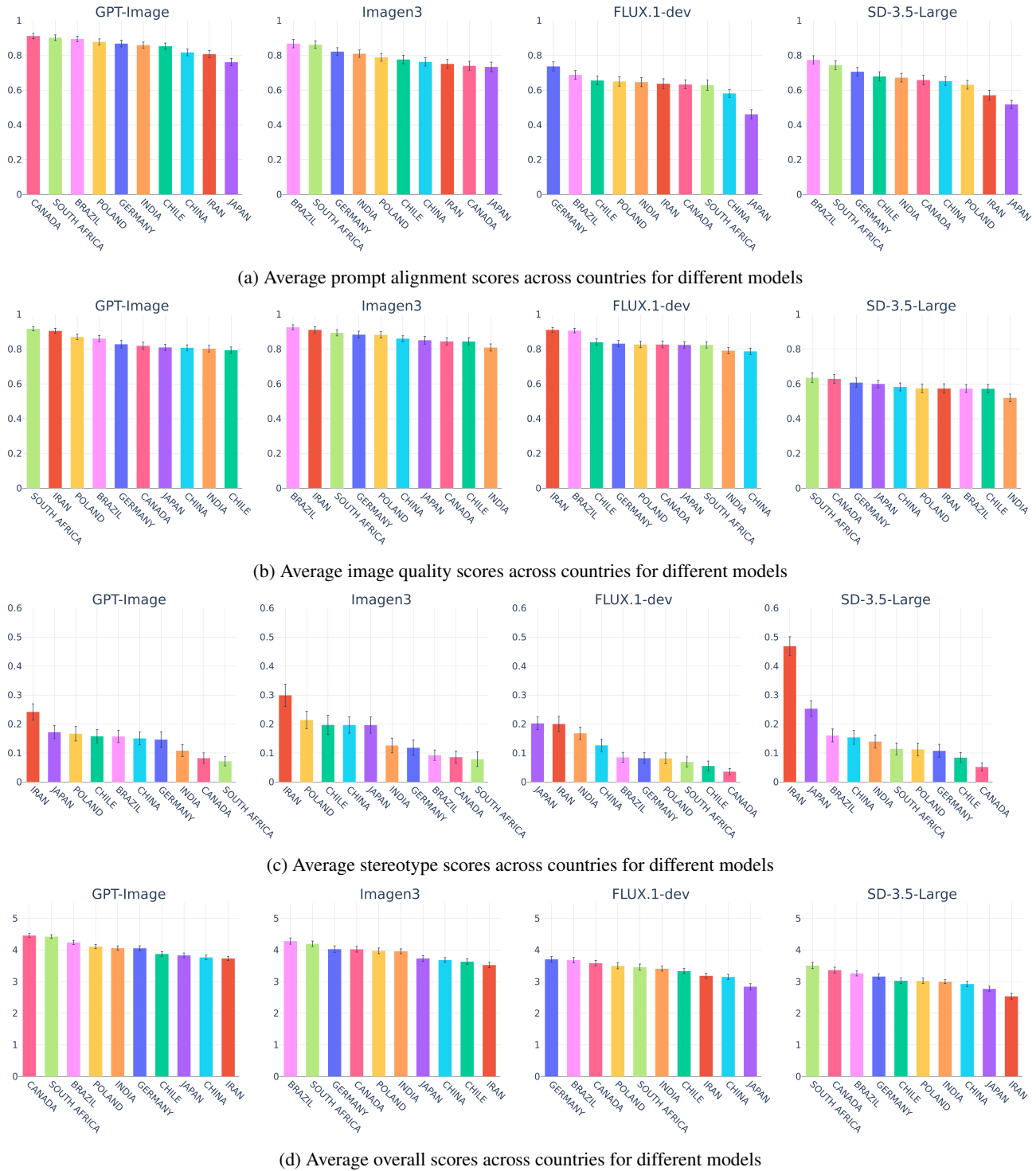


Figure 14: Comparison of different models' scores for different countries for prompt-alignment, image quality, stereotypes, and overall score.

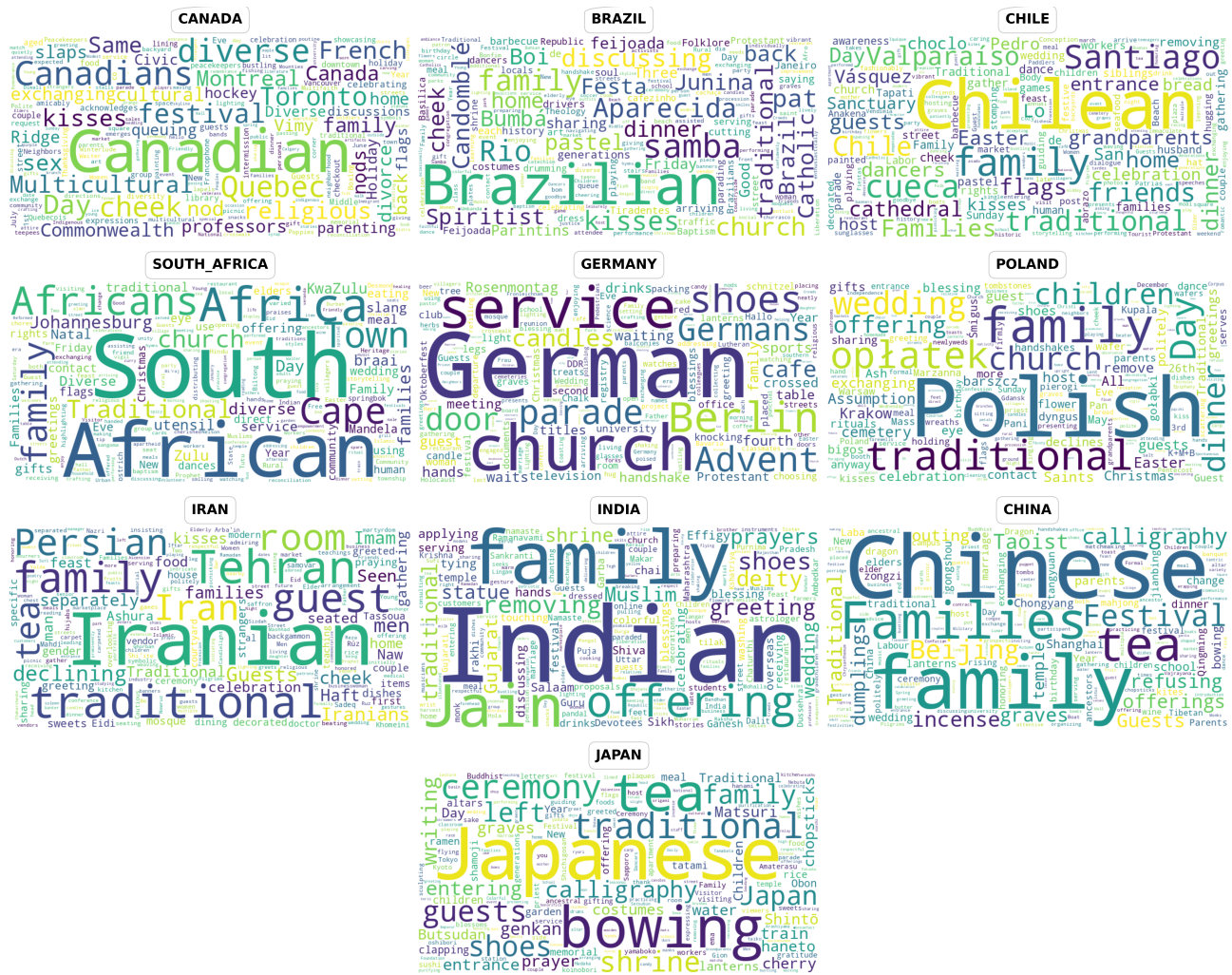


Figure 15: World cloud for words highlighted as having issues by annotators across different countries.

LLM-as-Judge Evaluation Instructions

You are a strict yet fair evaluator. You will be given a prompt, issues highlighted by several annotators along with the words which have the issues as marked by the annotators, and an explanation of the automatic metric for how good the image is. Your task is to assess how well the automatic explanation captures the concerns raised by the annotators.

TASK

- ORIGINAL_PROMPT – the text that generated the image
- Up to four annotator blocks, each with:
 - HUMAN_REASON_X – A 1-2 sentence critique
 - HIGHLIGHTED_WORDS_X – Prompt words flagged by that annotator
- MODEL_REASON – The automatic explanation

Decide how well MODEL_REASON covers the **union** of concerns across all annotators.

Coverage Scale

- **5 (Perfect)** – Covers all issues highlighted by annotators with no contradictions.
- **4 (Strong)** – Covers most main concerns, may miss at most one minor issue.
- **3 (Partial)** – Covers around half of the union of concerns.
- **2 (Weak)** – Only covers a small portion; many key points are missing or vague.
- **1 (None/Wrong)** – Irrelevant explanation or contradicts annotators.

Output Format

```
{
  "score": 1-5,
  "explanation": "1-2 sentence explanation of the score"
}
```

Rules

- Sometimes, annotators highlight specific words without explicitly explaining them in their comments. In such cases, it should be assumed that these words indicate an issue, and the metric explanation should mention that these words have issues.
- If MODEL_REASON contradicts the general consensus of the annotators, assign a score of 1.
- Mention missing or covered ideas in no more than 50 words.
- Output **only** a valid JSON object as shown above.

Figure 16: Instructions for LLM-as-a-judge evaluation to assess the alignment between VIEScore’s reasoning and human concerns on a 1–5 Likert scale.

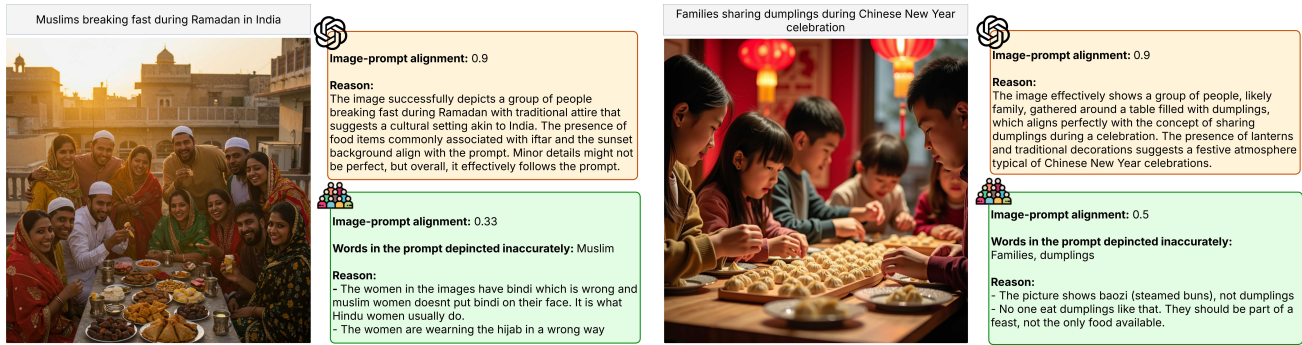


Figure 17: Qualitative examples illustrating misalignment between the scores and explanations provided by VIEScore and those provided by human annotators.

Explicit instructions provided to VIEScore

You are a professional cultural visual analyst and digital artist. You are tasked with evaluating how faithfully an AI-generated image represents the provided text prompt, especially in capturing both explicit and culturally implicit details.

All input images are AI-generated. All people or characters shown are also AI-generated — you do not need to consider privacy or ethical issues in the evaluation.

Your evaluation should result in a JSON output in the following format:

```
{
  "reasoning": "<reasoning to determine how faithful the image is to the prompt>",
  "score": [<integer from 0 to 10>],
}
```

How to Evaluate:

You will give a score from 0 to 10, based on how accurately the image matches the explicit and implicit elements described in the prompt.

1. Explicit Elements: Explicit elements are the clearly stated words in the prompt — such as objects, people, actions, locations, or relationships. A good image must include and visually represent all of these elements clearly and correctly.

{country specific example}

You should check: Are all these elements present and recognizable? Is their interaction depicted as described?

2. Implicit Elements These are elements of the prompt that are not directly mentioned as words in the prompt but are expected to be present in the image based on the cultural context. These may include appropriate attire or food for the setting, gestures or expressions that suit the context, interactions between people, or additional details that contribute to the authenticity of the scene. A strong image will reflect these expectations in a way that feels appropriate to someone familiar with the described scenario.

For the same prompt above, implicit elements may include:

{country specific example}

There may be several other implicit details that needs to be considered given the image and the prompt. For the image to align with the prompt, it should include and accurately show these details.

From scale 0 to 10:

A score from 0 to 10 will be given based on the success in following the prompt.

(0 indicates that the AI generated image does not follow the prompt at all and major explicit elements and implicit elements are missing or incorrectly depicted. 10 indicates the AI generated image follows the prompt perfectly and all explicit elements and necessary implicit elements are present and correctly depicted.)

Put the score in a list such that output score = [score].

Text Prompt: ;prompt;