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

Shravan Nayak<sup>1,2</sup> Mehar Bhatia<sup>1,3</sup> Xiaofeng Zhang<sup>1,2</sup>

Verena Rieser<sup>5</sup> Lisa Anne Hendricks<sup>5</sup> Sjoerd van Steenkiste<sup>4</sup>
Yash Goyal<sup>6</sup> Karolina Stańczak<sup>1,3,7</sup> Aishwarya Agrawal<sup>1,2</sup>

<sup>1</sup>Mila – Quebec AI Institute, <sup>2</sup>Université de Montréal, <sup>3</sup>McGill University, <sup>4</sup>Google Research, <sup>5</sup>Google DeepMind, <sup>6</sup>Samsung - SAIT AI Lab, Montreal, <sup>7</sup>ETH AI Center

Correspondence: shravan.nayak@mila.quebec

## **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 - where missed cues can stereotype communities and undermine usability. In this work, we present the first study to systematically quantify the alignment of T2I models and evaluation metrics with respect to both explicit (stated) as well as implicit (unstated, implied by the prompt's cultural context) 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 find that 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 show that existing T2I evaluation metrics correlate poorly with human judgments of cultural alignment, irrespective of their internal reasoning. Collectively, our findings expose critical gaps, provide a concrete testbed, and outline actionable directions for developing culturally informed T2I models and metrics that improve global usability.

https://culturalframes.github.io

# 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

both 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 and 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 while omitting implicit elements central to the scene, (such as a sacred fire or officiating priest). We refer to these two classes as explicit (based on the words in the prompt) and implicit (unstated but implied by the prompt's cultural context) expectations. 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 and Nushi, 2023).

This necessitates evaluation practices that not only verify faithfulness to the explicit expectations but also assess the inference and contextualization of implicit cultural expectations. However, current T2I evaluation methodologies predominantly focus on the former by assessing explicit promptimage consistency using automated metrics (Hu et al., 2023; Hessel et al., 2021; Ku et al., 2024a). 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

<sup>&</sup>lt;sup>1</sup>The 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.

safety (Lee et al., 2023), typically using generic or Western-centric prompts. Consequently, current evaluation methods and benchmarks lack adequate representation of culturally nuanced and expectation-rich scenarios critical to diverse cultural contexts.

In response, we present the first systematic study of cultural alignment in T2I models covering both explicit and implicit 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 annotators with relevant cultural backgrounds, who provide fine-grained judgments of the generated images across four criteria (i) image-prompt alignment, decomposed into explicit and implicit expectations; (ii) image quality; (iii) stereotype presence; and (iv) an overall score. For (i)-(iii), annotators

also provide explanations for their ratings. This scheme enables fine-grained analysis of T2I models' performance, providing rich insights. We find that state-of-the-art T2I models not only struggle with depicting implicit expectations but also clearly stated explicit ones. In fact, models fail to meet cultural expectations 44% of the time across countries. Among these instances, the failure rate for explicit expectations is unexpectedly high, averaging 68%, while the rate for implicit expectations is also substantial at 49%. We also observe that image quality varies by country, and stereotypes are flagged more often for Asian countries, particularly Japan and Iran, across all models.

Furthermore, we compare these human assessments with existing T2I evaluation metrics to demonstrate that current measures correlate poorly with human judgments of cultural alignment. In particular, VLM-based evaluators that produce rationales (e.g., VIEScore) give explanations that do not align with human reasons, calling into question the interpretability of their scores in culturally sensitive settings. Collectively, our findings lead to a discussion on actionable directions for developing more culturally informed T2I models and evaluation methodologies. These include turning our insights into better prompting strategies for models and metrics and, prospectively, using CULTURAL-FRAMES to align models and calibrate metrics.

Dataset	Countries	<b>Cultural Focus</b>	Prompts	Models	Annot.	Explicit Align.	Implicit Align.	Stereotype Flag	Explanation for Ratings	Human Eval. of Metrics
CUBE (Kannen et al., 2025)	8	Concept-centric	1,000	2	_	/	Х	Х	X	X
CultDiff (Bayramli et al., 2025)	10	Concept-centric	1,500	3	4,500	✓	X	×	X	✓
MC-SIGNS (Yerukola et al., 2025)	85	Gestures	288	2	1,408	X	X	✓	X	X
ViSAGe (Jha et al., 2024)	135	People	_	1	_	X	X	/	X	X
UCOGC (Zhang et al., 2024)	30	Material and non-material	752	3	67,620	✓	X	×	X	X
CulturalFrames (Ours)	10	Social practices & norms	983	4	10,000	✓	✓	✓	1	✓

Table 1: Comparison of cultural evaluation datasets for text-to-image generation, showing coverage (countries, cultural focus, scale), and whether the dataset supports evaluation of cultural alignment (explicit/implicit), stereotype flagging, and explanations for ratings. The final column shows if it includes human evaluation of 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 embeddingbased metrics, e.g., FID (Heusel et al., 2018), Inception Score (Salimans et al., 2016), and the textimage alignment via pretrained vision-language embeddings, e.g., CLIPScore (Hessel et al., 2021) and DinoScore (Ruiz et al., 2023). More recently, reward models trained on human preferences such as HPSv2 (Wu et al., 2023), 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 quality through questionanswering or reasoning, such as TIFA (Hu et al., 2023), DSG (Cho et al., 2024), V2QA (Yarom et al., 2023), VQAScore (Lin et al., 2025), UnifiedReward (Wang et al., 2025), DeQA (You et al., 2025), VIEScore (Ku et al., 2024b), and LLM-Score (Lu et al., 2023).

## **Cultural Alignment Evaluation of T2I models.**

T2I models struggle to accurately and respectfully represent cultural elements, leading to misrepresentation of cultural concepts and values (Ventura et al., 2025; Prabhakaran et al., 2022; Struppek et al., 2023). A growing body of work highlights various cultural biases, such as nationality-based biases (Jha et al., 2024; Alsudais, 2025), skin tone bias (Cho et al., 2023), social biases across gender, race, age, and geography (Bird et al., 2023;

Naik and Nushi, 2023). Other works focus on geographic representation (Basu et al., 2023; Hall et al., 2024), showing skewed generations towards Western contexts and study cultural adaptation through image editing (Khanuja et al., 2024).

Several recent benchmarks aim to probe cultural alignment in T2I systems (see Tab. 1). CUBE (Kannen et al., 2025) and CULTDIFF (Bayramli et al., 2025) focus on concept-centric cultural elements like food and landmarks across 8–10 countries and also do not assess implicit alignment or collect explanations for ratings. UCOGC (Zhang et al., 2024) covers 30 countries and evaluates both material and non-material culture, but does not address implicit cues, stereotype flagging, or human evaluation of metrics. MC-SIGNS (Yerukola et al., 2025) targets gestures, and VISAGe (Jha et al., 2024) focuses on portrayals of people, mainly emphasizing stereotype and offensiveness flags without assessing alignment or collecting explanations.

Qadri et al. (2025), a concurrent study, examines the limitations of evaluation practices through culturally grounded evaluations in three South Asian countries and advocates for "thick evaluations." Our work aligns with this emphasis but differs in being larger-scale and quantitative across countries, models, and metrics. As shown in Tab. 1, to the best of our knowledge, this is the first systematic quantification of how T2I models and metrics align with implicit cultural expectations.

## 3 CULTURALFRAMES

We detail our 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 knowl-

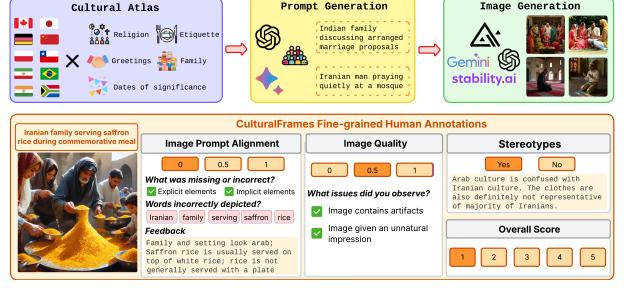


Figure 2: Overview of our data annotation pipeline. Cultural Atlas assertions are used to generate prompts using LLMs and human filtering (to-left, top-middle), which are converted to images using T2I models (top-right). Human annotators then evaluate alignment, quality, stereotypes, and overall score with detailed feedback (bottom).

edge, 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).<sup>2</sup>

# 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 Cultur-

alAtlas for the selected countries and their potential to induce prompts that elicit both explicit (*i.e.*, elements directly mentioned in the prompt) and implicit (*i.e.*, not mentioned in the prompt but inferred from shared cultural commonsense and needed for cultural authenticity) cultural expectations.

# 3.3 Data Generation Pipeline

Building on the 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 T2I models capture both explicit and implicit cultural expectations. Fig. 2 summarizes the data generation and human image annotation pipeline.

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-40 (OpenAI, 2024) using designed instructions (see § A.1). These instructions guide the model to embed cultural expectations into the prompts for realistic and observable everyday scenarios. Next, we use GPT-40 (OpenAI, 2024) and Gemini (Gemini Team, 2024) to automatically

<sup>&</sup>lt;sup>2</sup>We 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.

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 § A.2). Example assertions and prompts from our benchmark are shown in Tab. 3 in § A.1.

Image Generation. We generate images using four state-of-the-art T2I 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, active by default. To keep the evaluation practical and consistent across models, we generate one image per model per prompt. While this may appear limiting, our analysis (Appendix A.6) shows that output diversity across generations is generally low, and key issues identified by annotators tend to generalize across multiple outputs. In Fig. 15, 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 § B. To ensure high data quality, we filtered for attentive annotators and ensured a minimum of 25 unique, culturally knowledgeable workers<sup>3</sup> per country. We collect data from three annotators for each country using the Prolific<sup>4</sup> platform. Our annotation process captures detailed, multifaceted 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 relation 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. See Fig. 2 (bottom) for an example of human annotation for different criteria for an image-prompt pair.

# 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. 14 in § A.3.

**Images.** We generate images for our prompt set using both open- and closed-source models. While open-source models produce an image for every prompt, the safety filters of closed-source models block a subset of generations. This issue is most noticeable with Imagen3, which filters out 290 prompts—29.5% of the prompts, primarily due to policies against depicting children <sup>5</sup>. For comparison, GPT-40 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 three annotators. To measure agreement among raters, we compute Krippendorff's alpha (Krippendorff, 2013), obtaining 0.32 for prompt alignment, 0.28 for image quality, and 0.36 for the overall score. These scores are comparable to or better than prior works on cultural understanding in T2I models (Kannen et al., 2025; Bayramli et al., 2025). A detailed comparison with prior works, along with factors influencing the agreement scores, is provided in Appendix A.7.

<sup>&</sup>lt;sup>3</sup>Annotators were selected based on the following criteria: born in the country, a 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.

<sup>4</sup>https://www.prolific.com/

<sup>&</sup>lt;sup>5</sup>We requested an exemption from the provider to bypass these filters and will incorporate the missing images if access is granted

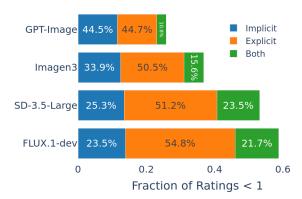


Figure 3: Distribution of image-prompt alignment errors (score <1) by model, grouped by error type: implicit, explicit, or both. Bars show fraction of total errors; percentages indicate each type's share per model.

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. 4 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. Our findings (Fig. 5 and Fig. 21) 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 (see Appendix C.2), showing a discernible trend where Asian countries tend to assign lower scores across multiple criteria.

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. 3 shows that GPT-Image has the lowest overall image-prompt alignment error rate (ratings < 1), with its errors roughly evenly split between implicit and explicit types. In contrast, other models, particularly SD-3.5-Large and FLUX, exhibit 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.

Which words do models most frequently misinterpret? Fig. 22 displays every word in the

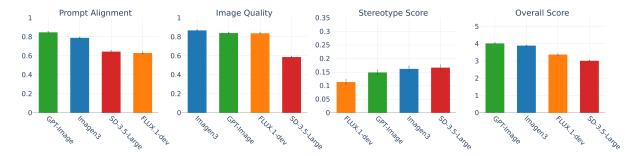


Figure 4: Human evaluation results for selected T2I models. From left to right: 1) Prompt Alignment (0–1, 1=perfect), 2) Image Quality (0–1, 1=highest), 3) Stereotype Score (0–1, 0=none), 4) Overall Score (1–5, 5=best).

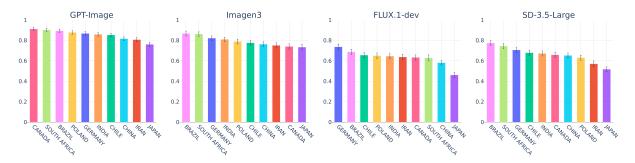


Figure 5: Prompt alignment scores across countries for a given model.

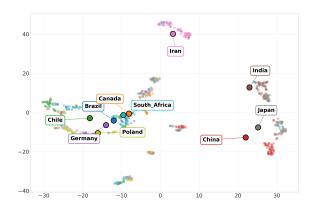


Figure 6: tSNE plot of Imagen3 images. Labeled markers show image embedding centroids per country.

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.

In what way do models fail 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 transformer<sup>6</sup> and cluster them using HDBScan (Campello et al., 2013). We then prompt GPT-40 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. American outputs suffered from poor regional specificity and inaccurate depictions of people's appearances. Similarly, 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. As shown in Fig. 6, images generated by Imagen3 for Asian countries form distinct clusters, while those from other

<sup>6</sup>https://huggingface.co/sentence-transformers/ all-mpnet-base-v2

 $<sup>^{7}</sup>$ https://huggingface.co/openai/clip-vit-large-patch14

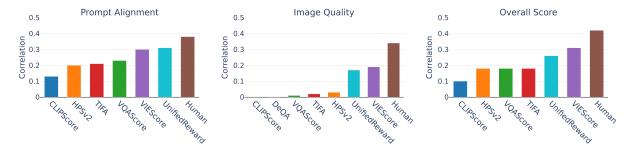


Figure 7: Spearman rank correlation of various T2I evaluation metrics with human ratings across three criteria: prompt alignment, image quality, and overall score. Human denotes the human-human Spearman rank correlation.

regions lack such clear grouping. This finding is corroborated by annotators in Europe and South America, who struggle to identify country-specific visual cues in generated images, indicating that the model fails to capture cultural distinctiveness.

# 6 Evaluating T2I Metrics on CULTURALFRAMES

Metrics analyzed. We analyze six representative metrics, each reflecting a different evaluation paradigm: CLIPScore (Hessel et al., 2021) is an embedding-based metric that computes cosine similarity between CLIP embeddings of the image and prompt. HPSv2 (Wu et al., 2023) enhances CLIPScore by fine-tuning the CLIP model on human preference data. TIFA (Hu et al., 2023) uses a VQA-based framework to assess faithfulness. We use GPT-4o-mini for question generation and Qwen2.5-VL-32B-Instruct (Team, 2025) as the answering model. VQAScore (Lin et al., 2025), UnifiedReward (Wang et al., 2025), and VI-EScore (Ku et al., 2024b) leverage vision-language models to evaluate generated images. For VQAScore, we leverage the CLIP-FlanT5 model introduced in the original VQAScore paper, use UnifiedReward-qwen-7B based on Qwen2.5-VL-7B for UnifiedReward, and use GPT-40 as VLM for VIEScore, which provides both a score and a textual reason for its assessment. Finally, we evaluate DeQA (You et al., 2025), a VLM trained specifically for image-quality assessment.

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. 7). UnifiedReward, an open-source reward model, slightly edges the best closed-model setup, VIEScore, on prompt alignment, achieving a Spearman correlation of 0.31 compared to 0.30 for the

latter. While this is below the human-human agreement of 0.38, it notably outperforms all other metrics. In contrast, TIFA exhibits a lower correlation, potentially because it only accounts for explicit elements mentioned in the prompt. This highlights a gap between metric design and 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 closely by UnifiedReward. The relatively stronger performance of HPSv2 may be attributed to its training on image pairs, with human preference likely driven by image quality, potentially making it more sensitive to visual appeal. By contrast, DeQA, despite being trained specifically for image-quality assessment on standard IQA datasets, shows near-zero correlation ( $\approx 0.0$ ) on our benchmark, likely due to domain and distribution shift between the data used to train DeQA and CULTURAL-FRAMES. Taken together, 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 (human-human: 0.42), with UnifiedReward close behind. Notably, HPSv2, despite being trained on human preferences, shows relatively poor performance, likely due to limited annotator and prompt diversity in the human preference dataset it was trained on. CLIPScore consistently underperforms, indicating limitations as a general-purpose evaluation metric, particularly for culturally sensitive image assessments. Overall, these results suggest that VLM-based metrics have the upper hand in capturing culturally grounded human.

Do explanations provided by VLM-based metrics capture the mistakes human raters

**highlight?** To further analyze the overall bestperforming metric on our benchmark, VIEScore, we evaluate whether its generated explanations reflect the issues raised by human annotators. We consider only cases where at least two annotators flagged mistakes with substantiated reasons. We adopt an LLM-as-a-judge setup, 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. 23. To mitigate potential model biases, we collect scores from 4 different LLMs - GPT-40, Gemini 2.5 Flash, Claude3.5-Sonnet, DeepSeek-Chat – and aggregate them per instance. To calibrate the LLM's judgments, we provided five in-context examples corresponding to varying quality levels. Additionally, we manually review 100 judge-provided scores, sampled across countries, and confirm they produce consistent, high-quality assessments. The results reveal that VIEScore's explanations achieve an average rating of 2.19/5 (std: 1.19), indicating only partial overlap with human rationale. This suggests that current metrics have substantial room to improve alignment with human judgments and reasoning. Some qualitative examples are provided in Tab. 9.

# 7 Discussion

Based on our analysis of cultural misalignment in text-to-image models and their evaluation metrics, we highlight three key directions for improvement.

Can culturally informed prompt expansion improve cultural alignment? CulturalFrames prompts are concise, leaving many cultural aspects implicit for the model to infer. For example, the prompt "a bride and groom exchanging vows at their Hindu wedding" omits elements like the priest or the sacred fire. To examine whether making these cues explicit can improve generations, we build on our analysis of model failures and expand the 20 lowest-scoring prompts per country using Gemini-2.5-Flash (see § C.1 for details). We generate images using Flux.1-Dev, and evaluate image-prompt alignment with VIEScore, the metric that best correlates with human judgments. Prompt expansion improves the average VIEScore from 7.3 to 8.4, demonstrating that culturally informed expansion helps models better capture details important to humans.

Can we improve metric performance through explicit instructions? Current T2I metrics are

not explicitly guided to account for implicit and explicit prompt elements. To test whether such guidance helps, we modify VIEScore by replacing GPT-40's instructions with novel annotation guidelines we developed for human raters (see Fig. 24). This raises correlation with human ratings from 0.30 to 0.32 (significant at 95% confidence) and improves explanation alignment with human rationales using the same LLM-as-a-judge setup from 2.19 to 2.37 on a 5-point scale. These results suggest that culturally informed instruction design can improve both scores and rationales. Nonetheless, the metric's reasoning still lags human rationale, highlighting the need for richer cultural knowledge and training beyond prompt design.

Does explicit training of VLMs to judge images improve culturally aligned evaluation? Current VLMs used for evaluation are typically not trained to judge images, raising the question of whether such training aids cultural alignment. We compare UnifiedReward (Wang et al., 2025), trained on diverse multimodal preference data to judge images, with its backbone model (Qwen2.5-VL-7B). Despite not targeting cultural scenarios, UnifiedReward achieves higher correlations with human judgments – image–prompt alignment (0.31 vs. 0.17), image quality (0.17 vs. 0.01), and overall score (0.28 vs. 0.14) – and even surpasses VI-EScore in alignment (0.31 vs. 0.30). This suggests preference-based judge training can significantly improve cultural alignment in metric scores.

## 8 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.

## 9 Limitations

Our study faces limitations due to our data collection methods and the scope of the CULTUR-ALFRAMES. 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 per 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.

# 10 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 foun-

dation for such future investigations.

# 11 Acknowledgements

We would like to thank Saba Ahmadi, Qian Yang, Ankur Sikarwar and Rohan Banerjee for their help with early pilots for prompt generation and image rating. We also thank the Mila IDT team for their technical support and for managing the computational resources. Additionally, Aishwarya Agrawal received support from the Canada CIFAR AI Chair award throughout this project. Karolina Stańczak was supported by the Mila P2v5 grant, the Mila-Samsung grant, and by an ETH AI Center post-doctoral fellowship. This project was generously funded by a research grant from Google.

# References

Muhammad Farid Adilazuarda, Sagnik Mukherjee, Pradhyumna Lavania, Siddhant Shivdutt Singh, Alham Fikri Aji, Jacki O'Neill, Ashutosh Modi, and Monojit Choudhury. 2024. Towards measuring and modeling "culture" in LLMs: A survey. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 15763–15784, Miami, Florida, USA. Association for Computational Linguistics.

Canfer Akbulut, Kevin Robinson, Maribeth Rauh, Isabela Albuquerque, Olivia Wiles, Laura Weidinger, Verena Rieser, Yana Hasson, Nahema Marchal, Iason Gabriel, William Isaac, and Lisa Anne Hendricks. 2025. Century: A framework and dataset for evaluating historical contextualisation of sensitive images. In *International Conference on Learning Representations (ICLR)*.

Abdulkareem Alsudais. 2025. Analyzing how text-to-image models represent nationalities in everyday tasks. *Preprint*, arXiv:2504.06313.

Abhipsa Basu, R. Venkatesh Babu, and Danish Pruthi. 2023. Inspecting the geographical representativeness of images from text-to-image models. In 2023 IEEE/CVF International Conference on Computer Vision (ICCV), pages 5113–5124.

Zahra Bayramli, Ayhan Suleymanzade, Na Min An, Huzama Ahmad, Eunsu Kim, Junyeong Park, James Thorne, and Alice Oh. 2025. Diffusion models through a global lens: Are they culturally inclusive? *Preprint*, arXiv:2502.08914.

Charlotte Bird, Eddie L. Ungless, and Atoosa Kasirzadeh. 2023. Typology of risks of generative text-to-image models. *Preprint*, arXiv:2307.05543.

Ricardo JGB Campello, Davoud Moulavi, and Jörg Sander. 2013. Density-based clustering based on

- hierarchical density estimates. In *Pacific-Asia conference on knowledge discovery and data mining*, pages 160–172. Springer.
- Jaemin Cho, Yushi Hu, Roopal Garg, Peter Anderson, Ranjay Krishna, Jason Baldridge, Mohit Bansal, Jordi Pont-Tuset, and Su Wang. 2024. Davidsonian scene graph: Improving reliability in fine-grained evaluation for text-to-image generation. *Preprint*, arXiv:2310.18235.
- Jaemin Cho, Abhay Zala, and Mohit Bansal. 2023. Dall-eval: Probing the reasoning skills and social biases of text-to-image generation models. *Preprint*, arXiv:2202.04053.
- Nassim Dehouche and Kullathida Dehouche. 2023. What's in a text-to-image prompt? the potential of stable diffusion in visual arts education. *Heliyon*, 9(6):e16757.
- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, Kyle Lacey, Alex Goodwin, Yannik Marek, and Robin Rombach. 2024. Scaling rectified flow transformers for high-resolution image synthesis. *Preprint*, arXiv:2403.03206.
- Dan Friedman and Adji Bousso Dieng. 2023. The vendi score: A diversity evaluation metric for machine learning. *Preprint*, arXiv:2210.02410.
- Gemini Team. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *Preprint*, arXiv:2403.05530.
- Christian Haerpfer, Ronald Inglehart, Alejandro Moreno, Christian Welzel, Kseniya Kizilova, Marta Lagos, Juan Diez-Medrano, Pippa Norris, Eduard Ponarin, and Bi Puranen. 2022. World Values Survey: Round seven country-pooled datafile version 3.0. Madrid, Spain & Vienna, Austria: JD Systems Institute & WVSA Secretariat.
- Melissa Hall, Candace Ross, Adina Williams, Nicolas Carion, Michal Drozdzal, and Adriana Romero Soriano. 2024. Dig in: Evaluating disparities in image generations with indicators for geographic diversity. *Preprint*, arXiv:2308.06198.
- Stuart Hall. 1980. Encoding/decoding. In Stuart Hall, Dorothy Hobson, Andrew Lowe, and Paul Willis, editors, *Culture, Media, Language: Working Papers in Cultural Studies*, pages 63–87. Hutchinson, London.
- Jochen Hartmann, Yannick Exner, and Samuel Domdey. 2025. The power of generative marketing: Can generative AI create superhuman visual marketing content? *International Journal of Research in Marketing*, 42(1):13–31.
- 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 crosscultural NLP. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6997–7013, Dublin, Ireland. Association for Computational Linguistics.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 2021. CLIPScore: A reference-free evaluation metric for image captioning. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 7514–7528, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2018. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Preprint*, arXiv:1706.08500.
- Geert Hofstede, Gert Jan Hofstede, and Michael Minkov. 2010. *Cultures and organizations: software of the mind: intercultural cooperation and its importance for survival*, 3rd edition. McGraw-Hill, New York; London.
- Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang, Mari Ostendorf, Ranjay Krishna, and Noah A Smith. 2023. TIFA: Accurate and interpretable text-toimage faithfulness evaluation with question answering. *Preprint*, arXiv:2303.11897.
- Kaiyi Huang, Chengqi Duan, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. 2025. T2I-CompBench++: An Enhanced and Comprehensive Benchmark for Compositional Text-to-Image Generation. *IEEE Transactions on Pattern Analysis Machine Intelligence*, (01):1–17.
- Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. 2023. T2I-CompBench: A comprehensive benchmark for open-world compositional text-to-image generation. *Advances in Neural Information Processing Systems*, 36:78723–78747.
- Imagen-Team-Google. 2024. Imagen 3. *Preprint*, arXiv:2408.07009.
- Akshita Jha, Vinodkumar Prabhakaran, Remi Denton, Sarah Laszlo, Shachi Dave, Rida Qadri, Chandan Reddy, and Sunipa Dev. 2024. ViSAGe: A global-scale analysis of visual stereotypes in text-to-image generation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12333–12347, Bangkok, Thailand. Association for Computational Linguistics.
- Dongfu Jiang, Max Ku, Tianle Li, Yuansheng Ni, Shizhuo Sun, Rongqi Fan, and Wenhu Chen. 2024. GenAI arena: An open evaluation platform for generative models. *Preprint*, arXiv:2406.04485.

- Nithish Kannen, Arif Ahmad, Marco Andreetto, Vinodkumar Prabhakaran, Utsav Prabhu, Adji Bousso Dieng, Pushpak Bhattacharyya, and Shachi Dave. 2025. Beyond aesthetics: Cultural competence in text-to-image models. *Preprint*, arXiv:2407.06863.
- Simran Khanuja, Sathyanarayanan Ramamoorthy, Yueqi Song, and Graham Neubig. 2024. 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*, pages 10258–10279, Miami, Florida, USA. Association for Computational Linguistics.
- Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. 2023. Pick-apic: An open dataset of user preferences for text-to-image generation. *Preprint*, arXiv:2305.01569.
- K. Krippendorff. 2013. *Content Analysis: An Introduction to Its Methodology*. SAGE Publications.
- Max Ku, Dongfu Jiang, Cong Wei, Xiang Yue, and Wenhu Chen. 2024a. VIEScore: Towards explainable metrics for conditional image synthesis evaluation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12268–12290, Bangkok, Thailand. Association for Computational Linguistics.
- Max Ku, Dongfu Jiang, Cong Wei, Xiang Yue, and Wenhu Chen. 2024b. VIEScore: Towards explainable metrics for conditional image synthesis evaluation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12268–12290, Bangkok, Thailand. Association for Computational Linguistics.
- Max Ku, Tianle Li, Kai Zhang, Yujie Lu, Xingyu Fu, Wenwen Zhuang, and Wenhu Chen. 2024c. Imagen-Hub: Standardizing the evaluation of conditional image generation models. *Preprint*, arXiv:2310.01596.
- Black Forest Labs. 2024. Flux. https://github.com/black-forest-labs/flux.
- Tony Lee, Michihiro Yasunaga, Chenlin Meng, Yifan Mai, Joon Sung Park, Agrim Gupta, Yunzhi Zhang, Deepak Narayanan, Hannah Teufel, Marco Bellagente, Minguk Kang, Taesung Park, Jure Leskovec, Jun-Yan Zhu, Fei-Fei Li, Jiajun Wu, Stefano Ermon, and Percy S Liang. 2023. Holistic evaluation of texto-image models. In *Advances in Neural Information Processing Systems*, volume 36, pages 69981–70011. Curran Associates, Inc.
- 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 *Computer Vision ECCV 2024*, pages 366–384, Cham. Springer Nature Switzerland.
- Soumaya Loukili, Lotfi Elaachak, and Abdelhadi Fennan. 2025. Finetuning stable diffusion models for

- email marketing text-to-image generation. In *Innovations in Smart Cities Applications Volume 8*, pages 524–535, Cham. Springer Nature Switzerland.
- Yujie Lu, Xianjun Yang, Xiujun Li, Xin Eric Wang, and William Yang Wang. 2023. LLMScore: Unveiling the power of large language models in text-to-image synthesis evaluation. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Adyasha Maharana, Darryl Hannan, and Mohit Bansal. 2022. StoryDALL-E: Adapting pretrained text-to-image transformers for story continuation. In *Computer Vision ECCV 2022*, pages 70–87, Cham. Springer Nature Switzerland.
- Marshall McLuhan. 1966. *Understanding Media: The Extensions of Man.* Signet Books, New York.
- Mosaica. 2024. The cultural atlas. https://cultural atlas.sbs.com.au/.
- Ranjita Naik and Besmira Nushi. 2023. Social biases through the text-to-image generation lens. *Preprint*, arXiv:2304.06034.
- OpenAI. 2024. GPT-40 system card. *Preprint*, arXiv:2410.21276.
- OpenAI. 2025. Introducing 40 image generation. http s://openai.com/index/introducing-40-image-generation/.
- Vinodkumar Prabhakaran, Rida Qadri, and Ben Hutchinson. 2022. Cultural incongruencies in artificial intelligence. *Preprint*, arXiv:2211.13069.
- Rida Qadri, Mark Diaz, Ding Wang, and Michael Madaio. 2025. The case for "thick evaluations" of cultural representation in AI. *Preprint*, arXiv:2503.19075.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, and 1 others. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmJ.R.
- Charvi Rastogi, Tian Huey Teh, Pushkar Mishra, Roma Patel, Zoe Ashwood, Aida Mostafazadeh Davani, Mark Diaz, Michela Paganini, Alicia Parrish, Ding Wang, Vinodkumar Prabhakaran, Lora Aroyo, and Verena Rieser. 2024. Insights on disagreement patterns in multimodal safety perception across diverse rater groups. *Preprint*, arXiv:2410.17032.
- Charvi Rastogi, Tian Huey Teh, Pushkar Mishra, Roma Patel, Ding Wang, Mark Díaz, Alicia Parrish, Aida Mostafazadeh Davani, Zoe Ashwood, Michela Paganini, Vinodkumar Prabhakaran, Verena Rieser, and Lora Aroyo. 2025. Whose view of safety? a deep dive dataset for pluralistic alignment of text-to-image models. *Preprint*, arXiv:2507.13383.

- Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. 2023. DreamBooth: Fine tuning text-to-image diffusion models for subject-driven generation. *Preprint*, arXiv:2208.12242.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. 2022. Photorealistic text-to-image diffusion models with deep language understanding. *Preprint*, arXiv:2205.11487.
- Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. 2016. Improved techniques for training gans. *Preprint*, arXiv:1606.03498.
- Lukas Struppek, Dom Hintersdorf, Felix Friedrich, Manuel Br, Patrick Schramowski, and Kristian Kersting. 2023. Exploiting cultural biases via homoglyphs in text-to-image synthesis. *Journal of Artificial Intelligence Research*, 78:1017–1068.
- Qwen Team. 2025. Qwen2.5-vl.
- Mor Ventura, Eyal Ben-David, Anna Korhonen, and Roi Reichart. 2025. Navigating cultural chasms: Exploring and unlocking the cultural POV of text-to-image models. *Transactions of the Association for Computational Linguistics*, 13:142–166.
- Yixin Wan, Arjun Subramonian, Anaelia Ovalle, Zongyu Lin, Ashima Suvarna, Christina Chance, Hritik Bansal, Rebecca Pattichis, and Kai-Wei Chang. 2024. Survey of bias in text-to-image generation: Definition, evaluation, and mitigation. *Preprint*, arXiv:2404.01030.
- Yibin Wang, Yuhang Zang, Hao Li, Cheng Jin, and Jiaqi Wang. 2025. Unified reward model for multimodal understanding and generation. *Preprint*, arXiv:2503.05236.
- Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. 2023. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis. *Preprint*, arXiv:2306.09341.
- Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. 2023. Imagereward: learning and evaluating human preferences for text-to-image generation. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, pages 15903–15935.
- Michal Yarom, Yonatan Bitton, Soravit Changpinyo, Roee Aharoni, Jonathan Herzig, Oran Lang, Eran Ofek, and Idan Szpektor. 2023. What you see is what you read? improving text-image alignment evaluation. In *Advances in Neural Information Processing Systems*, volume 36, pages 1601–1619. Curran Associates, Inc.

- Akhila Yerukola, Saadia Gabriel, Nanyun Peng, and Maarten Sap. 2025. Mind the gesture: Evaluating AI sensitivity to culturally offensive non-verbal gestures. *Preprint*, arXiv:2502.17710.
- Zhiyuan You, Xin Cai, Jinjin Gu, Tianfan Xue, and Chao Dong. 2025. Teaching large language models to regress accurate image quality scores using score distribution. In *IEEE Conference on Computer Vision and Pattern Recognition*.
- Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Ben Hutchinson, Wei Han, Zarana Parekh, Xin Li, Han Zhang, Jason Baldridge, and Yonghui Wu. 2022. Scaling autoregressive models for content-rich texto-image generation. *Preprint*, arXiv:2206.10789.
- Lili Zhang, Xi Liao, Zaijia Yang, Baihang Gao, Chunjie Wang, Qiuling Yang, and Deshun Li. 2024. Partiality and misconception: Investigating cultural representativeness in text-to-image models. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, CHI '24, New York, NY, USA. Association for Computing Machinery.

# **Appendix**

# **Table of Contents**

	Page
A. CULTURALFRAMES	
A.1 Prompt Generation	14
A.2 Prompt Filtering	15
A.3 Prompt Distribution	15
A.4 Image Generation	15
A.5 Prompt-Image Examples	15
A.6 Image Generation Analysis	15
A.7 Inter Human Agreement	16
B. Image Rating	
B.1 Rating Interface	18
B.2 Annotator Demographics	18
C. Text-to-Image Models' Analysis	
C.1 Prompt Expansion Case Study .	19
C.2 Image Quality Analysis	20
C.3 Model Ranking Analysis	20
C.4 Model Scores Analysis	20
C.5 Word Cloud of Issues	33
D. Text-to-Image Metrics' Analysis	
D.1 LLM-as-a-Judge Protocol	34
D.2 VIEScore Failure Examples	35
D.3 VIEScore Revised Instructions	36

## A 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** 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 select 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. Example assertions and prompts from our benchmark are shown in Tab. 3.

We began by generating a small number of prompts per category, which were evaluated by

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.

Assertion (CulturalAtlas)	<b>Generated Prompts</b>
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.	<ul><li>(1) Grandchildren touching grandfather's feet at an Indian temple.</li><li>(2) Indian village elder blessing children during harvest festival.</li></ul>
Religion (Iran): Most Iranians believe in Islam, but due to politicization, many younger citizens have withdrawn. Devout followers often practice privately at home.	<ul><li>(1) Iranian family praying together at home.</li><li>(2) Elderly Iranian man praying in a quiet mosque.</li></ul>

Table 3: Examples of assertions in CulturalAtlas for two categories *greetings* in India and *religion* in Iran and corresponding generated prompts.

human annotators to assess whether the scenarios were both visually depictable and culturally appropriate (see Section A.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-40 (OpenAI, 2024) used across different categories are provided in Figures 9, 10, 11, 12, 13.

# A.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. 8, 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.3 Prompt Distribution Across Categories**

Fig. 14 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.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.5 Prompt-Image Examples

Some examples of prompts along with images generated using different models are provided in Fig. 15.

## A.6 Image Generation Analysis

We generate only one image per prompt due to the practical constraints of our annotation budget



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?



Figure 8: 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.

and the need to maintain a manageable scale. Despite this limitation, we believe our findings remain meaningful and generalizable, particularly given the known low diversity in model outputs (Kannen et al., 2025). To substantiate this, we conducted two additional analyses:

Quantifying Image Diversity for Cultural-Frames We analyze the diversity of generated images using the best-performing open-source model, Flux.1-Dev (Labs, 2024). For every prompt in CulturalFrames, we generate 4 images using different random seeds. We then embed these images using the CLIP model (ViT-L/14@336px) (Radford et al., 2021) and compute the Vendi Score (Friedman and Dieng, 2023), which reflects the effective number of distinct images in a set. Across all prompts, we find an average Vendi score of 1.5 (standard deviation 0.3) for 4 images, indicating that only 1.5 unique images are produced out of 4 on average. This result confirms the low diversity previously reported in the literature.

Checking Generalization of Annotator Comments To assess whether annotator observations generalize to other images, we manually inspect 4 images each for 20 prompts from India, Poland, and China, countries whose cultural norms our authors are familiar with. These prompts were selected because annotators had already identified cultural issues in the single-image setup.

In all 20 cases, at least three out of four images exhibited the same cultural issues that had been pre-

viously flagged. This finding strongly reinforces our initial observations and demonstrates that these issues generalize consistently across multiple generations. Tab. 4 provides qualitative examples of prompts and the cultural issues highlighted by annotators, along with the number of images in which these issues were observed.

These results support our claim that even with multiple generations, the same cultural issues tend to persist. This is likely due to the limited diversity of current models. Therefore, while we only use one image per prompt in our main evaluation, our findings do generalize to multi-image settings for current generation systems. Lastly, we believe that the rich explanations collected from annotators can be extremely valuable for future work that studies model biases in multi-image generation settings.

## A.7 Inter Human Agreement

To establish that our inter-annotator agreement is well within the field's norms, we quantitatively compare our country-level Krippendorff's Alpha and Fleiss' Kappa scores against published values from two closest benchmarks, CUBE (Kannen et al., 2025) and CultDiff (Bayramli et al., 2025). For Krippendorff's Alpha, across both image-prompt alignment and image-quality, CULTURALFRAMES's country-level scores consistently match and often exceed the lower bounds of CUBE's reported ranges (e.g., CUBE's image-prompt alignment: 0.09–0.58 vs. CULTURALFRAMES: 0.24–0.42). Similarly, for Fleiss' Kappa,

Prompt	Observed Cultural Issue	Prevalence
Visitors removing shoes before entering a Hindu temple in India.	Annotators commented that people were not removing their shoes, and many were still wearing shoes as they entered the temple.	4/4
Chinese couple receiving parental blessings in traditional attire.	Annotators observed that there were no parents visible in the images.	4/4
Families sharing dumplings during Chinese New Year celebration.	Annotators complained that the food shown in the image is "baozi" rather than dumplings.	4/4
Children float Marzanna doll down Polish river to end winter.	Annotators complained that there is no Marzanna doll in the image.	4/4
Families cooking rice dishes under festive decorations during Pongal.	Annotators pointed out that there was a fire over the rice kept in the dish.	3/4

Table 4: Examples of Persistent Cultural Issues Across Multiple Image Generations

Gender	Iran	Chile	Germany	Japan	India	China	Canada	South Africa	Brazil	Poland	Average
Male	0.68	0.68	0.80	0.60	0.80	0.70	0.73	0.84	0.82	0.74	0.74
Female	0.74	0.80	0.82	0.53	0.73	0.60	0.80	0.77	0.84	0.72	0.72

Table 5: Average image-prompt alignment scores by gender and country. The numbers highlighted have a difference greater than 0.5.

Age Group	Germany	Iran	Chile	Japan	India	China	Canada	South Africa	Brazil	Poland	Average
18-24	0.84	0.71	0.77	0.69	0.74	0.65	0.75	0.80	0.83	0.76	0.75
25-44	0.78	0.67	0.78	0.61	0.78	0.71	0.73	0.77	0.85	0.72	0.74
45+	0.76	0.71	0.45	0.57	0.67	0.73	0.76	0.78	0.77	0.72	0.68

Table 6: Average image-prompt alignment scores by age groups and country. The numbers highlighted have a difference greater than 0.5.

our agreement on prompt alignment (0.179–0.406) and image quality (0.157–0.341) is noticeably higher than CultDiff's general figures (0.07–0.17). For the overall score, where both datasets share a 1–5 scale, our agreement (0.06–0.14) is comparable. Importantly, CulturalFrames attains these agreement levels despite requiring raters to judge more subtle, implicit cultural cues than the more object-level signals in the benchmarks. We credit this strong performance to our meticulously designed evaluation framework, where we iteratively update instructions and filter workers to ensure high data quality. To understand inter-human agreement for CulturalFrames better, we quantitatively and qualitatively analyze several key factors:

**Do people of different genders rate images differently?** For every country, we split the annotations by gender and calculate the mean scores provided by each gender for the image-prompt alignment criteria. Our data is predominantly annotated by peo-

ple who identify as *male* or *female*, except Japan, where 1 annotator did not identify with either gender. Hence, we present the analysis across only these two categories of gender. We make sure to include only those prompt-image instances (2248 of them) where we have ratings from both genders to ensure fair evaluation.

Tab. 5 provides the average image-prompt alignment scores provided by male and female annotators. We begin by examining the overall average scores across gender groups: males score 0.74 and females score 0.72, resulting in a modest gap of 0.02. This difference is slightly higher than the 0.01 gap observed when annotations are randomly split, suggesting that gender may play a minor but measurable role in rating variation. However, this effect appears more pronounced when analyzed at the country level.

Several countries in Tab. 5 exhibit notable gender-based differences in cultural alignment

scores. Chile shows the largest gap, with females scoring 0.80 and males 0.68. China also reflects a considerable difference, with males scoring significantly higher, 0.70, and females 0.60. Canada, India, Japan, and South Africa also demonstrate moderate differences, with females and males differing by over 0.06. These gaps may reflect differences in perception, interpretation, or cultural sensitivity across genders in line with previous works that study gender based variations in T2I evaluation (Rastogi et al., 2024, 2025). Despite these variations, some countries like Germany, Brazil, and Poland show more consistent scores between male and female annotators.

Do people from different age groups rate images differently? For each country, we categorize annotators into three age groups (18-24, 25-44, 45+). This corresponds to Gen Z, GenX, and millennials, respectively. We make sure to include only those prompt-image instances (2407 of them) where we have ratings from two of the three age groups (as ensuring all three age groups annotated an instance filtered a lot of annotations, as we collect only 3 human annotations for a prompt-image pair) to ensure fair evaluation. We calculate the average prompt alignment scores and report them in Tab. 6.

The age-wise analysis reveals clear generational differences in how cultural alignment is rated. On average, annotators aged 18–24 give the highest scores (0.75), followed closely by the 25–44 group (0.74), while the 45+ group gives notably lower scores (0.68). This 0.07 drop between the youngest and oldest age groups is substantially higher than the 0.01 difference observed when annotations are randomly split (3-way random split, each pairwise difference was 0.01) and differences are calculated, suggesting that age meaningfully influences evaluation behavior.

On a country level, annotators aged 18–24 assign the highest scores the most number of times (5/10 countries), followed by the 24-44 age group (4/10 countries), suggesting they may be more optimistic, lenient, or culturally flexible. This trend is most prominent in Chile (0.77 for 18-24 vs. 0.45 for 45+), Japan (0.69 for 18-24 vs. 0.57 for 45+), and India (0.78 for 24-44 vs. 0.67 for 45+). In contrast, older participants (45+) tend to give lower scores, indicating more critical assessments, possibly due to deeper cultural anchoring. Countries like Iran, South Africa, and Canada exhibit relatively stable scores across age groups, suggesting less genera-

tional variance in perception. This analysis underscores the importance of considering age-based diversity when evaluating subjective alignment tasks, as perspectives can shift meaningfully across generations.

# Are people's sensitivities to the same issues dif-

ferent? We analyze whether annotators may provide similar reasoning for their judgments but assign different alignment scores, indicating varying sensitivities to the same issue. We observe such instances in our dataset and argue that this variation is not annotator noise, but a natural outcome of subjective interpretation in value-centric evaluations. The rationales we collect alongside each score are critical in making sense of these differences, offering insight into annotators' thought processes and allowing us to study the nuances behind disagreement, rather than dismissing them as inconsistencies. We include qualitative examples below to illustrate this phenomenon in Tab. 7.

# Do people flag different issues for the same im-

age? We observe that in a small number of cases, different annotators identify different issues in the same image, which can stem from their diverse cultural backgrounds and lived experiences. What one annotator flags as a misrepresentation may not even register to another, highlighting the subjectivity inherent to cultural evaluation, which could result in different scores. We provide qualitative examples to illustrate this phenomenon in Tab. 8. Further, we note that the combination of diverse perspectives provided by the annotators in these cases collectively covers a broad spectrum of potential issues, leading to a more holistic and robust understanding of cultural expectations.

# **B** Image Rating

# **B.1** Rating Interface

We develop a custom interface for collecting image ratings. Fig. 16 and Fig. 17 show the detailed instructions we provide to the annotators for rating images. Fig. 18 shows the interface where annotators rate images.

# **B.2** Annotator Demographics

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

Prompt	Annotator 1 Comment	Score	Annotator 2 Comment	Score
Chinese villagers gathering for Laba Festival porridge feast	It is not Laba-style por- ridge	0	That is not Laba porridge. Laba porridge contains at least 8 ingredients.	0.5
Casual hug between German friends at a cafe	I cannot see a hug. Beer is not typical for a cafe.	0	The hug is hard to see. The scene doesn't look casual. The person is drinking beer in a cafe.	0.5
Polish family passing pierogi platter during afternoon meal	These are not pierogi.	0	These are not traditional Polish dumplings	0.5
Canadians lining up outside mall for Boxing Day sales	There isn't much of a line, nor does this suggest the people here are Canadi- an/in Canada	0.5	There is no explicit depiction of Canadians, nor is there a line present, nor is the "line" outside the mall.	0

Table 7: Qualitative examples of different sensitivities in scores shown by annotators for the same or similar issues. The score to the right of the annotator comment is the rating the annotator provided.

Prompt	Annotator 1 Comment	Annotator 2 Comment
Sikh children learning Gurbani in an Indian classroom	Children look more like a for- eigner than an Indian sikh.	Girl Sikh children don't wear turbans.
Families sharing fish meal on Good Friday in Brazil	The image does not depict a fish meal like Brazilian people eat it. In fact, the fish looks raw, so it's weird.	I can't see a Brazilian family in this photo; there is nothing that indicates it. It looks more like Asian people.
Traditional African ceremony in KwaZulu Natal province	Men aren't wearing the traditional dress, which would include animal hide. This is an important part of Zulu culture and wouldn't be changed.	There is nothing resembling KwaZulu Natal province, including the clothing and the scenery.

Table 8: Qualitative examples of different annotators providing different reasons for their ratings.

# C Text-to-Image Models' Analysis

# C.1 Prompt Expansion Case Study

Building on the insights gathered from our detailed analysis of model failures, we propose a simple but effective prompt expansion strategy. Our annotator rationales revealed recurring patterns in what models tend to overlook, such as missing cultural objects, family members, inaccuracies in settings, and mood. To test whether explicitly including these overlooked details in the prompt improves generation authenticity, we selected the 20 lowest-scoring prompts from each country (200 prompts in

total across 10 countries) and expanded the prompts using an LLM (Gemini-2.5-Flash). The LLM was given the instructions detailed in Fig. 19.

We generate images using the Flux model (the best open-source model) for these expanded prompts, and use VIEScore to measure the image-prompt alignment accuracy. We use VIEScore as it is the metric that correlates the most with human judgments. We see that there is a consistent improvement of VIEScore (overall score) from 7.3 to 8.4 upon prompt expansion, indicating that careful prompt expansion could indeed help in model

improvement.

# **C.2** Image Quality Analysis

We analyze images that receive low ratings and find that low quality scores are most often linked to artifacts (70.4% of cases) and an unnatural impression (50.9% on average). Across models, SD-3.5-Large accounts for the largest share of both artifact flags (54.4%) and 'unnatural' flags (43.2%). Notably, Flux-1.0-dev and GPT-Image also show high 'unnatural' shares ( $\approx$ 24% and  $\approx$ 22%, respectively). Our qualitative analysis indicates that 'unnatural' is typically triggered by global coherence issues where scenes or cultural elements seem implausible for the setting, whereas 'artifacts' reflects local distortions (e.g., blur, distortions).

## **C.3** Model Ranking Analysis

We try to understand if there is a preferred model across countries and how they rank. For prompt alignment (see Fig. 20), 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. 21, 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.

# **C.4** Model Scores Analysis

Fig. 21 presents country-level scores across the four criteria evaluated by human annotators.

## **Prompt Instructions (Greeting)**

#### **Purpose**:

We want to test whether text-to-image models can accurately capture a country's distinct greeting practices. You will be given:

- 1. A country name
- 2. A short excerpt on greeting norms: an implicit description of how people in this country typically greet each other, or some information that relates to greeting customs.

#### 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 greeting values and norms. 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 greeting norm or value mentioned in the excerpt can be observed. 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. This includes avoiding verbal greetings that cannot be depicted 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}

Figure 9: Instructions used to generate prompts for the greeting category

# **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}

Figure 10: Instructions used to generate prompts for the religion category

#### **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 "arrving 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:**

 $\{ \verb|incontext_examples_negative| \}$ 

Figure 11: Instructions used to generate prompts for the etiquette category

## **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}

Figure 12: Instructions used to generate prompts for the family category

## **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.

#### Vour 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}

Figure 13: Instructions used to generate prompts for the dates of significance category



Figure 14: Distribution of prompts from different categories across countries.

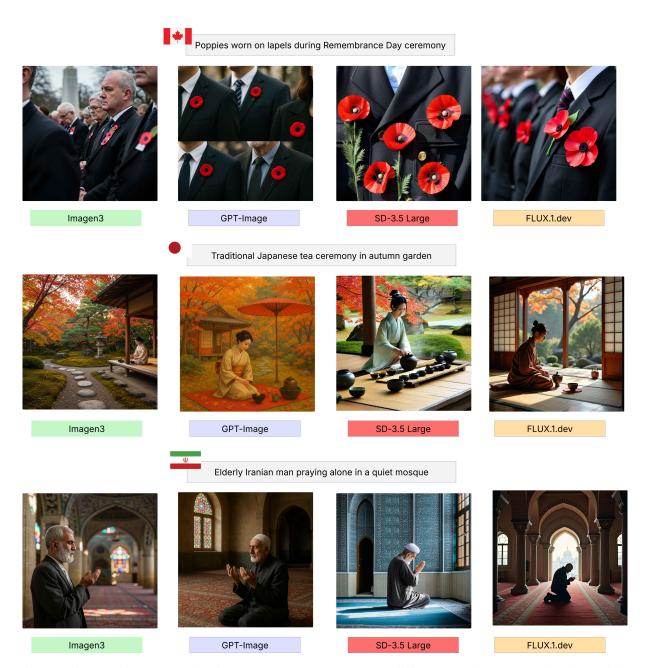


Figure 15: Prompt-image examples from CulturalFrames across different countries generated by the models.

## **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 10 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.
  - $\,\circ\,$  You  $\pmb{\mathsf{must}}$  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.

Figure 16: Prompt alignment instructions provided to the annotators. The example shown varies depending on the countries.

#### 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  ${\it Yes}$ , 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 unsual sense that make objects or the scene hard to recognize or understand.

Note: You can access this information anytime during rating by clicking the 1 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 17: Instructions given to annotators for stereotype, image quality, and overall score criteria.

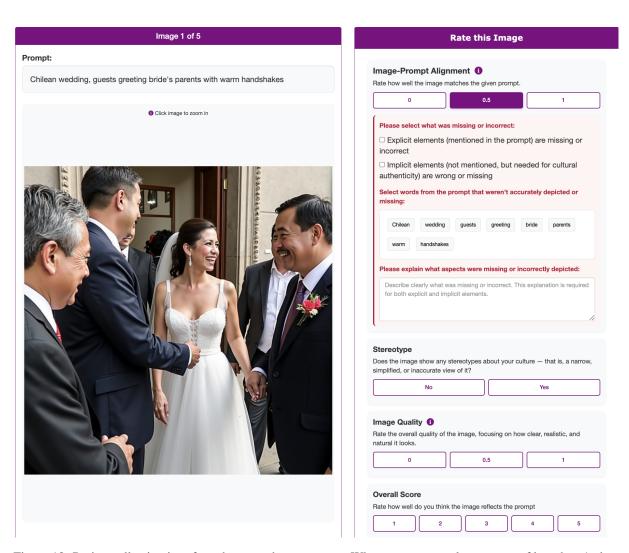


Figure 18: Rating collection interface shown to the annotators. When annotators select a score of less than 1, they need to give detailed feedback regarding explicit and implicit expectations, along with selecting the problematic words.

## **Prompt Expansion Instructions**

#### **Purpose:**

You are an expert in cultural nuance and creative image generation. Your task is to expand the following brief cultural prompt into a more detailed and descriptive one suitable for a state-of-the-art AI image generator. The goal is to create a visually rich and culturally authentic image.

Original Prompt: "original\_prompt"

Instructions for Expansion: Enrich the prompt by adding vivid details across these categories:

- Setting and Environment: Describe the specific location, time of day, lighting, and background elements.
- People and Demographics: Detail the family members' approximate ages, their relationships to one another, their attire, and their expressions.
- Objects and Food: Specify the types of food on the table, the serving dishes, and any other relevant objects in the scene.
- Cultural Atmosphere and Mood: Capture the overall feeling of the scene—is it lively, warm, formal, or relaxed?
- Artistic Style: Suggest a photographic style (e.g., "cinematic, warm lighting, shallow depth of field, 35mm film look").

Combine all of the above details into a single, cohesive, descriptive paragraph.

The output should be in the following format:

Expanded: <expanded\_prompt>

Figure 19: Instructions provided to an LLM to generate expanded prompts.



Figure 20: Model ranking across countries for different criteria (1 is the highest rank). Countries are grouped by geographical proximity.

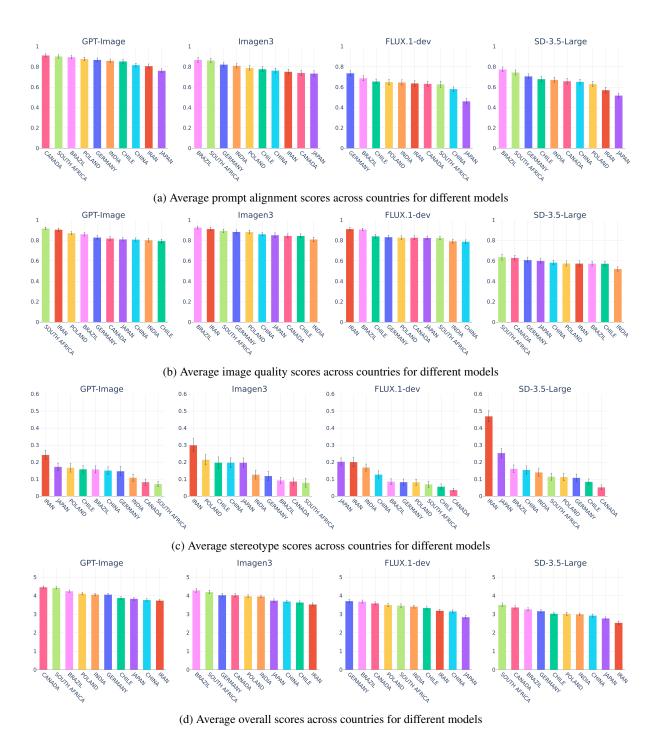


Figure 21: Comparison of different models' scores for different countries for prompt-alignment, image quality, stereotypes, and overall score.

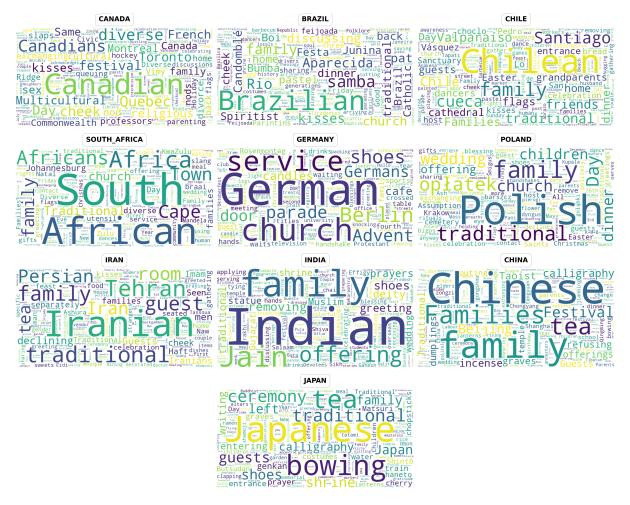


Figure 22: World cloud for words highlighted as having issues by annotators across different countries.

#### **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 23: Instructions for LLM-as-a-judge evaluation to assess the alignment between VIEScore's reasoning and human concerns on a 1–5 Likert scale.



**Prompt**: Japanese family celebrating New Year's at home with osechi-ryori

**Human rationale**: Usually, the Japanese don't make their face white on New Year's. The way of holding chopsticks is strange. The foods look like suhi rather than osechi ryori.

**VIEScore rationale**: The image accurately reflects the prompt, showing a Japanese family in traditional kimonos celebrating New Year. Osechi-ryori is prominently featured, and the home setting with sliding doors matches cultural expectations. Implicit cues like traditional attire, joyful expressions, and a family-oriented atmosphere further enhance authenticity.



**Prompt**: Devotees lighting candles at a Jain temple during Paryushana in India.

**Human rationale**: Paryushana would not be engraved inside the temple. There are no candles, and how can they light them with their hands?

**VIEScore rationale**: The image effectively depicts devotees lighting candles in what appears to be a Jain temple setting. The architecture and attire of the individuals are consistent with an Indian cultural setting, and the script in the background corresponds to the theme of Paryushana. The ambiance and lighting convey a religious and serene atmosphere. The missing point could be a more explicit reference to Paryushana elements.



**Prompt**: Faithful lighting candles at Aparecida shrine in Brazil

**Human rationale**: The church shown in the generated image doesn't look like the Aparecida Shrine but a generic catholic church.

VIEScore rationale: The image captures a large number of people in a religious setting, with many candles lit, fitting the description of faithful lighting candles. The architecture resembles a shrine, and the stained glass windows add to the spiritual ambiance. The scene closely aligns with the Aparecida shrine in Brazil. The only minor point could be about distinct visual cues specific to Brazil, which might not be strongly evident.

Table 9: Examples showing mismatches between human and VIEScore rationales.

#### **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 need 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:

Figure 24: Updated instructions provided to VIEScore, similar to those human raters use to judge images.