

---

# Exposing Blindspots: Cultural Bias Evaluation in Generative Image Models

---

Huichan Seo<sup>1\*</sup>    Sieun Choi<sup>2\*</sup>    Minki Hong<sup>2\*</sup>    Yi Zhou<sup>1</sup>    Junseo Kim<sup>3</sup>  
Lukman Ismaila<sup>4</sup>    Naome Etori<sup>5</sup>    Mehul Agarwal<sup>1</sup>    Zhixuan Liu<sup>1</sup>    Jihie Kim<sup>2</sup>  
Jean Oh<sup>1</sup>

<sup>1</sup>Carnegie Mellon University, Pittsburgh, United States

<sup>2</sup>Dongguk University, Seoul, South Korea

<sup>3</sup>Delft University of Technology, Delft, Netherlands

<sup>4</sup>Johns Hopkins University, School of Medicine, Baltimore, United States

<sup>5</sup>University of Minnesota–Twin Cities, Minneapolis, United States

Corresponding author: Jean Oh(jeanoh@cmu.edu)

## Abstract

Generative image models produce striking visuals yet often misrepresent culture. Prior work has probed cultural dimensions primarily in text-to-image (T2I) systems, leaving image-to-image (I2I) editors largely underexamined. We close this gap with a unified, reproducible evaluation spanning six countries, an 8-category/36-subcategory schema, and era-aware prompts, auditing both T2I generation and I2I editing under a standardized, reproducible protocol that yields comparable model-level diagnostics. Using open models with fixed configurations, we derive comparable diagnostics across countries, eras, and categories for both T2I and I2I. Our evaluation combines standard automatic measures, a culture-aware metric that integrates retrieval-augmented VQA with curated knowledge, and expert human judgments collected on a web platform from country-native reviewers. To enable downstream analyses without re-running compute-intensive pipelines, we release the complete image corpus from both studies alongside prompts and settings. Our study reveals three recurring findings. First, under country-agnostic prompts, models default to Global-North, modern-leaning depictions and flatten cross-country distinctions, reducing separability between culturally distinct neighbors despite fixed schema and era controls. Second, iterative I2I editing erodes cultural fidelity even when conventional metrics remain flat or improve; by contrast, expert ratings and our culture-aware metric both register this degradation. Third, I2I models tend to apply superficial cues (palette shifts, generic props) rather than context- and era-consistent changes, frequently retaining source identity for Global-South targets and drifting toward non-photorealistic styles; attribute-addition trials further expose weak text rendering and brittle handling of fine, culture-specific details. Taken together, these results indicate that culture-sensitive edits remain unreliable in current systems. By standardizing prompts, settings, metrics, and human evaluation protocols—and releasing all images and configurations—we offer a reproducible, culture-centered pipeline for diagnosing and tracking progress in generative image research. Project page: <https://seochan99.github.io/ECB/>

---

\*Equal contribution.



Figure 1: Representative cultural biases in T2I generations across six countries. Examples include Chinese–Japanese aesthetic conflation, mis-styled Indian weddings, Kenyan wildlife stereotypes, Korean attire misidentification, Nigerian safari mislocalization, and U.S. cultural miscues in food and religious ritual. Images are from FLUX.1 [schnell] fp8 and HiDream-I1-Dev.

## 1 Introduction

Text-to-image (T2I) and image-to-image (I2I) generative models have advanced rapidly in photorealism and controllability [1, 2, 3]. Yet T2I models exhibit systematic cultural bias rooted in imbalanced, skewed training data, where some regions and communities are over-represented while others are under-represented [4, 5]. This yields stereotyped or inaccurate portrayals of underrepresented groups, as shown in Figure 1. I2I editing is often used as a practical remedy by inserting or adjusting culture-specific elements [6], but errors persist across regions and frequently resurface over multiple edits [7, 8]. As a result, the burden of cultural correction shifts to users. We therefore focus on the *I2I editing loop*: a sequence of edits intended to align an image with a target culture, used to examine whether bias persists, attenuates, or reappears across iterations.

Turning to evaluation, cultural fidelity remains challenging to measure. Distributional or general-alignment metrics such as Fréchet Inception Distance (FID) [9] and CLIPScore [10] do not capture culture-specific attributes. Emerging cultural benchmarks [8, 11] help, but often rely heavily on human studies or use only author-generated synthetic images, limiting generalization [12].

We analyze five representative model families (Table 2) by first constructing a T2I base image set spanning six countries: China, India, Kenya, Korea, Nigeria, and the United States (U.S.). The set is balanced across eight cultural categories (Architecture, Art, Events, Fashion, Food, Landscape, People, Wildlife), further stratified into 36 subcategories (Table 1). Our prompts are *era-aware*, comprising *traditional*, *modern*, and *era-agnostic* variants, which allows us to probe temporal sensitivity in cultural understanding. To the best of our knowledge, this is the first study to systematically evaluate era-aware cultural competence in image models. Building on the base image set, we evaluate I2I cultural adaptation using three complementary experiments: (1) Multi-Loop Edit: five sequential edits that test whether the model preserves context while progressively improving cultural consistency; (2) Attribute Addition: stepwise insertion of five distinct culture-specific attributes to quantify cumulative competence and interference; and (3) Cross-Country Restylization: coherent transfer from a source country to a target country to assess adaptability and generalization.

For evaluation, we adopt automatic metrics—CLIPScore [10], DreamSim [13], and Aesthetic Score [14]—complemented by culture-centered assessments. We compare automatic metrics against two complementary evaluations: a VQA-based culture-aware metric and a human evaluation. This approach allows us to analyze the agreement, gaps, and limitations between automatic evaluations and human perception.

Beyond culture, our occupational audit shows persistent gender and skin-tone skews under neutral prompts, indicating embedded demographic priors. Three findings stand out: (i) under country-agnostic prompts, generations collapse to a United States–like, modern aesthetic; (ii) in iterative image-to-image editing, culture-relevant cues decay even as standard automatic metrics stay flat or improve, while a culture-aware assessment tracks human judgments; and (iii) I2I models often rely

on shortcut cues (palette, emblematic symbols) and often leave identity attributes (e.g., skin tone and facial traits) unchanged when restylizing to Global-South targets. Together, the findings expose limits of current metrics and benchmarks and argue for culture- and category-aware evaluation, alongside stronger data and training signals such as balanced curation and explicit debiasing/regularization.

Our contributions are as follows:

1. **Experimental design.** A geographically balanced, multi-domain schema (6 countries; 8 categories, 36 subcategories) and 3 complementary protocols (Multi-Loop Edit, Attribute Addition, Cross-Country Restylization) that contrast inherent T2I bias with I2I editing capability, including an era-aware prompt design to assess temporal awareness (an underexplored evaluation setting).
2. **Dataset release.** Public release of the complete image set, together with prompts and model/execution configurations, enabling downstream cultural analyses without re-running generation/editing or model reconfiguration.
3. **Evaluation framework.** An open-source evaluation platform that pairs automated metrics with culture-aware assessments and human studies, enabling triangulation across automated, culture-aware, and human judgments; we benchmark and validate metric behavior against human results to surface agreements and discrepancies.

## 2 Related Work

### 2.1 Generative Image Models

Generative image models rely on two closely related paradigms: T2I generation and I2I editing. Latent diffusion models [15, 2] set the T2I standard for fidelity and controllability, yet often miss fine-grained cultural fidelity. In response, I2I frameworks increase edit precision via context-aware conditioning, multi-step guidance, and cross-attention manipulation (e.g., FLUX.1 Kontext [16], HiDream-I1 [17], Qwen-Image-Edit [18], NextStep-1 [19]). These methods preserve locality and enable targeted attributes but are typically optimized for realism and generic prompt compliance rather than culture-specific correctness. Prior studies [20, 21, 22] show that stronger editing control alone does not mitigate entrenched biases: models can encode or even amplify stereotypes when data and objectives are not culture-aware. We, therefore, compare T2I baselines and their I2I editing loops to quantify initial cultural bias and the bias-correction burden placed on users.

### 2.2 Text-Image Datasets

The foundation of modern generative image models rests on large-scale text-image datasets. Early efforts like MS-COCO [23] and ImageNet [24] were predominantly Western-centric. Subsequent expansions, including LAION-5B [25], achieved unprecedented scale but inherently inherited severe cultural and geographical biases, notably English dominance and uneven global representation. While datasets like Dollar Street [26] and WIT [27] attempted to target diversity, their scope was limited by inherent editorial constraints. More recent benchmarks, including CCUB [28] and SCoFT [29], have focused on directly assessing cultural fidelity. However, these resources primarily prioritize concept coverage over the representation of nuanced cultural practices, ceremonies, or daily life contexts [30]. This persistent gap—where existing datasets remain ill-suited for assessing authentic cultural representation—motivates the creation of our targeted evaluation framework, which is the direct focus of this work.

### 2.3 Evaluation Metrics in Cultural Image Generation

The evolution of generative image models mandates corresponding advancements in output evaluation. Traditional metrics mainly assess image quality (e.g., FID [9], DreamSim [13], Aesthetic Score [14]) and text-image alignment (e.g., CLIPScore [10] and DinoScore [31]). Recent improvements leveraging human preference-trained models [32] and Large Language Model (LLM)-based metrics [33]. However, these approaches primarily capture realism and faithfulness to prompts, while neglecting cultural and contextual fidelity [8, 34]. Recent benchmarks such as CUBE [8], CULTDIFF [11], and CulturalFrames [30] aim to evaluate cultural alignment and bias, yet often depend on human judgments or limited internal datasets, constraining generalizability. To overcome

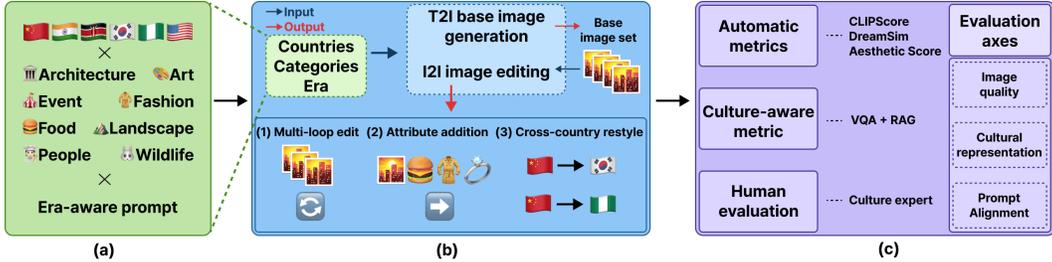


Figure 2: Overall framework overview. (a) Schema inputs: six countries, eight categories, and three era-aware prompts. (b) Experimental pipeline: T2I base generation and three I2I editing studies. (c) Multi-layered evaluation: integrating automatic, culture-aware metrics, and human evaluation.

this gap, we propose a comprehensive framework for evaluating authentic cultural representation in general generative outputs.

### 3 Experiment Design for Cultural Bias Evaluation

#### 3.1 Experimental Setup

We design two complementary studies: a T2I study for bias via image generation, and an I2I study for competence via editing culturally salient details. We evaluate across six geographically diverse countries (China, India, Kenya, Korea, Nigeria, U.S.), including Global South contexts to address underrepresentation in training data [30]. The overall framework is illustrated in Figure 2.

**Category schema.** Building on prior works [8, 35], our schema defines eight top-level categories (36 subcategories) covering both tangible artifacts (e.g., architecture, cuisine) and intangible practices (e.g., rituals, festivals), roles, and settings. This granular system probes cultural competence by enabling era-aware prompting and controlled cross-country comparisons to evaluate breadth and depth, rather than aesthetics alone systematically. The full schema is shown in Table 1.

**Models.** We conducted T2I and I2I experiments using state-of-the-art open-source models. Models were selected to capture capability diversity and ensure reproducibility. The full list of models is summarized in Table 2 with all settings detailed in the Appendix A.

**Evaluation metrics.** We evaluate models using two complementary dimensions: general-purpose and culture-aware. For general-purpose, we adopt CLIPScore, DreamSim, and Aesthetic Score to measure semantic alignment, cumulative edit distance, and visual quality. For culture-aware evaluation, we extend VQA- and RAG-based approaches [35, 36] by retrieving Wikipedia context via FAISS [37] and generating yes/no questions and answer with Qwen2.5-0.5B and Qwen2-VL-7B [38, 39]. Human ratings serve as the primary reference, enabling analysis of consistency and divergence between automatic and culture-aware metrics.

**Prompt design.** To disentangle temporal stylization from cultural identity, we use three prompt modes for each country/category/subcategory:

Table 1: Category schema used in our experiments.

category	subcategories	category	subcategories
<b>Architecture</b>	House, Landmark	<b>Landscape</b>	City, Countryside, Nature
<b>Art</b>	Dance, Painting, Sculpture	<b>Fashion</b>	Accessories, Clothing, Makeup
<b>Food</b>	Beverage, Dessert, Main dish, Snack, Staple food	<b>Wildlife</b>	Animal, Plant
<b>People</b>	Daily life, Athlete, Bride and groom, President, Soldier, Student, Teacher		

- Traditional: “Traditional {subcategory} in {Country}, photorealistic.”
- Modern: “Modern {subcategory} in {Country}, photorealistic.”
- Era-agnostic: “{Subcategory} in {Country}, photorealistic.”

This era-aware prompting allows us to examine how cultural depictions vary across specified and unspecified eras, revealing era-sensitive biases.

### 3.2 Text-to-Image (T2I) Experiment

We construct a standardized base image set to expose model-internal cultural priors under controlled text-only prompting. For each country and subcategory, we issue era-aware prompts (traditional, modern, and era-agnostic) while fixing all sampling parameters within each model family (Table 2, T2I column). This yields a comparable collection of generations across all settings, which is subsequently used for distributional analyses, traditional-modern scoring, and seeding the I2I editing experiments in Section 3.3. Fixing prompts and sampling settings isolates differences attributable to the models rather than to prompt phrasing or parameter drift.

### 3.3 Image-to-Image (I2I) Experiment

We evaluate cultural editing competence with three experiment designs using open-source I2I models—(Table 2, I2I column)—reflecting contemporary editing paradigms while remaining reproducible.

#### 3.3.1 Multi-loop edit: cultural consistency under iterative edits

To minimize cross-model domain shifts, base images are drawn from a single T2I model family. We then apply the instruction “Change the image to represent {era} {subcategory} in {Country}” for five successive rounds. This iterative experiment tests stability and path dependence in cultural editing—whether iterations converge toward culturally faithful attributes or drift by amplifying stereotypes.

#### 3.3.2 Attribute addition: compositional cultural grounding

To isolate composition effects, we begin with a neutral canvas: a genderless green mannequin on a white background devoid of cultural cues. For each country/model, a fixed five-step sequence controls (1) background, (2) local-script text rendering, (3) food, (4) clothing, and (5) traditional accessories. This sequence allows us to assess, within a single experiment, the model’s understanding across multiple facets commonly flagged as challenging in cultural benchmarks. Full stepwise prompts are provided in the Appendix F.1.

#### 3.3.3 Cross-country restyle: style transfer across countries

We use T2I base images from HiDream-I1-Dev (selected for strong T2I fidelity) and test whether I2I models can restyle an entire scene into a target country’s aesthetic while preserving the subject’s identity and pose. We employ a single minimal prompt—“Transform this image into the {CountryAdj} style.”—and report results under this condition; in preliminary trials, adding era/context/attribute cues did not materially improve cultural plausibility, so we do not vary prompt

Table 2: Models used in T2I and I2I experiments.

Family	T2I	I2I
Stable Diffusion [2]	Stable Diffusion 3.5 Medium	Stable Diffusion 3.5 Medium
FLUX.1 [16]	FLUX.1 [schnell] fp8	FLUX.1 Kontext [dev]
HiDream [17]	HiDream-I1-Dev	HiDream-E1.1
Qwen-Image [18]	Qwen-Image	Qwen-Image-Edit
NextStep [19]	NextStep-1-Large	NextStep-1-Large-Edit



Figure 3: Comparative samples for U.S. (top) and country-agnostic (bottom) prompts across two models. Left image: FLUX.1 [schnell] fp8; right image: HiDream-I1-Dev. Within each model, panels show (from left) *Bride and groom*, *Chef*, and *Farmer*. The close correspondence between rows illustrates the US-like default of country-agnostic prompts.

specificity further. This setup supports (i) qualitative audits of cultural appropriateness and context/era consistency and (ii) inspection of edit stability across multi-loop steps.

### 3.4 Expert Human Evaluation

We complement automatic metrics with an expert-group human evaluation, run on our web platform (ECB Human Survey) under a unified protocol. Raters evaluate only their own country (emic expertise), ensuring cultural authenticity in judgments. For each prompt, four candidates (step base/1/3/5) are displayed side-by-side. Raters assign two scores: 1-5 Likert scores for *Image Quality* and *Cultural Representation*. These two components are averaged to form the Human Quality Score (HQS). Raters also select *Best/Worst*, following recent practices in cultural assessment and editing evaluation [12, 40]. Image sets span eight categories across five models (Table 2). Each rater completes six tasks: five I2I-loop evaluation and one attribute addition evaluation. Expert raters require insider knowledge and language proficiency. Operational details, including the distribution of our expert raters, are provided in Appendix B.

## 4 Results

### 4.1 Text-to-Image (T2I)

We evaluate five T2I models with a unified CLIP-embedding backbone. Each generated image is embedded once and tagged by country (“country-agnostic” means no explicit tag). For clustering, we apply model-wise principal component analysis followed by  $k$ -means with  $K_m$  clusters, yielding a cluster index  $k(i)$  that exposes model-specific visual modes. The analysis proceeds along two axes: distributional proximity between countries and traditional–modern leaning within semantic categories.

**Cluster-proportion vectors.** For each model and country, we summarize the allocation of that country’s samples across latent modes as a probability vector.

$$\mathbf{p}_c^{(m)} = (p_{c,0}^{(m)}, \dots, p_{c,K_m-1}^{(m)}), \quad p_{c,k}^{(m)} = \frac{\#\{i : c(i) = c, k(i) = k\}}{\#\{i : c(i) = c\}}, \quad \sum_k p_{c,k}^{(m)} = 1. \quad (1)$$

This vector is a normalized mixture over recurring visual patterns; countries are close when their mass concentrates on the same modes.

**Distributional proximity.** We compare countries as distributions over modes by combining cosine similarity (directional alignment) with the Jensen–Shannon divergence (JSD; information overlap). The JSD is defined as

$$\text{JSD}(p, q) = \frac{1}{2}\text{KL}(p||m) + \frac{1}{2}\text{KL}(q||m), \quad m = \frac{1}{2}(p + q). \quad (2)$$

The per-model proximity that jointly rewards alignment and overlap is

$$h_{ab}^{(m)} = \frac{2 \cos(\mathbf{p}_a^{(m)}, \mathbf{p}_b^{(m)}) [1 - \text{JSD}(\mathbf{p}_a^{(m)}, \mathbf{p}_b^{(m)})]}{\cos(\mathbf{p}_a^{(m)}, \mathbf{p}_b^{(m)}) + [1 - \text{JSD}(\mathbf{p}_a^{(m)}, \mathbf{p}_b^{(m)})]}. \quad (3)$$

To reduce model idiosyncrasies, we report the fixed-effect average across models:

$$\bar{h}_{ab} = \frac{1}{M} \sum_{m=1}^M h_{ab}^{(m)}. \quad (4)$$

**Traditional–modern leaning.** Within each semantic category, we contrast images against category-specific prototypes to obtain a scale-invariant, signed leaning score (positive = traditional; negative = modern).

$$s_i = \cos(x_i, \mu_{\text{trad}}(g(i))) - \cos(x_i, \mu_{\text{mod}}(g(i))). \quad (5)$$

Country-level aggregates and uncertainty are

$$\bar{s}_c = \frac{1}{n_c} \sum_{i \in \mathcal{I}_c} s_i, \quad \text{SE}(\bar{s}_c) = \frac{\text{sd}(\{s_i : i \in \mathcal{I}_c\})}{\sqrt{n_c}}. \quad (6)$$

Category conditioning controls for topic mix; cosine focuses on directional similarity in embedding space; the difference-of-cosines removes norm sensitivity and yields a comparable, unitless score.

#### 4.1.1 Distributional Proximity and Cultural Defaults

For each model and country, we compute cluster-proportion vectors using Eq. 1, and measure pairwise country similarity using the proximity in Eq. 3, which combines cosine similarity with JSD. The model-averaged proximity is reported using Eq. 4. Across all five models, the strongest proximity is United States  $\leftrightarrow$  country-agnostic (mean 0.892, 95% CI [0.844, 0.931]), followed by China  $\leftrightarrow$  Korea (0.888, [0.835, 0.932]) and Kenya  $\leftrightarrow$  Nigeria (0.864, [0.811, 0.911]). A meta-analysis of between-model heterogeneity indicates tau-squared approximately zero overall, with a small but non-negligible value for Kenya–Nigeria (about  $8.5 \times 10^{-5}$ ), suggesting stable effects across architectures and training data. These results support a robust United States–like default for country-agnostic prompts and reveal regional clustering in East Asia and Sub-Saharan Africa. Model-level nearest-neighbor tallies likewise identify China–Korea and Kenya–Nigeria as mutual neighbors, while India frequently aligns with United States/country-agnostic, indicating partial assimilation. Figure 3 visually corroborates this pattern: United States and country-agnostic prompts yield nearly indistinguishable samples, consistent with a United States–like default.

#### 4.1.2 Traditional–Modern Leaning under Country-agnostic Prompts

Image-level traditional–modern scores are computed per model using Eq. 5, aggregated to country means with standard errors using Eq. 6, and evaluated for significance via within-category permutation tests. To address multiplicity across countries, we control the expected proportion of false positives among declared discoveries using the Benjamini–Hochberg false discovery rate (FDR) procedure at level 0.05. The cross-country dispersion of country means is significant in every model (four models:  $p=0.001$ ; Qwen-Image:  $p=0.002$ ), indicating systematic differences under country-agnostic prompts. The dominant pattern is a consistent modern lean for United States and country-agnostic across all five models (negative country means, per-model  $p$  values at most 0.006, FDR-significant). Other effects are smaller and model-dependent: for example, Kenya is traditional-leaning only in FLUX.1 [schnell] fp8 ( $p=0.001$ , FDR-significant), whereas India, Nigeria, Korea, and China are mixed or not significant in most settings (typically  $q$  values greater than 0.1). Category composition largely explains the aggregates. Clothing and sport tend to produce modern (negative) scores across models and countries—especially for United States and country-agnostic—whereas house, religious ritual, and landmark generally produce traditional (positive) scores for most countries but often flip to modern for United States and country-agnostic. These patterns indicate that the observed modern default arises from systematic topic/style composition rather than sampling noise. Per-country  $p$ -values and FDR-adjusted  $q$ -values are reported in Appendix C, Table 4.

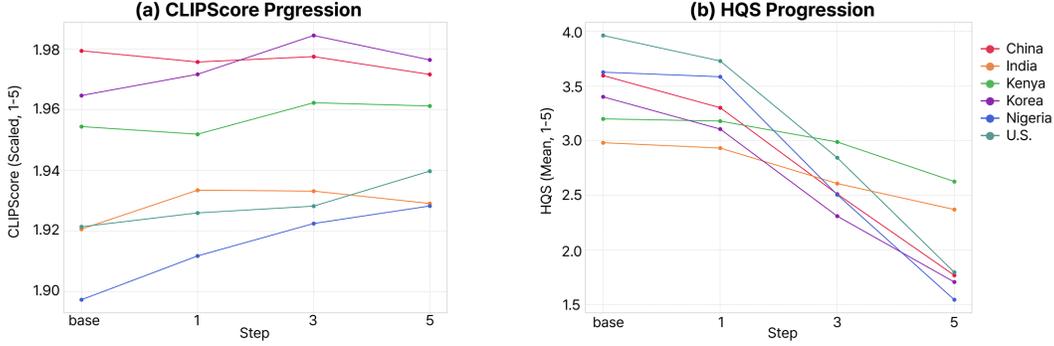


Figure 4: Divergence between Automated and Human Judgment in Iterative Editing. (a) CLIPScore trajectories remain largely stable or modestly increase from the base step to step 5. (b) Human Quality Score (HQS) sharply declines across all countries. This pronounced divergence highlights the failure of traditional automatic metrics to track the perceptible cultural degradation that human raters consistently penalize.

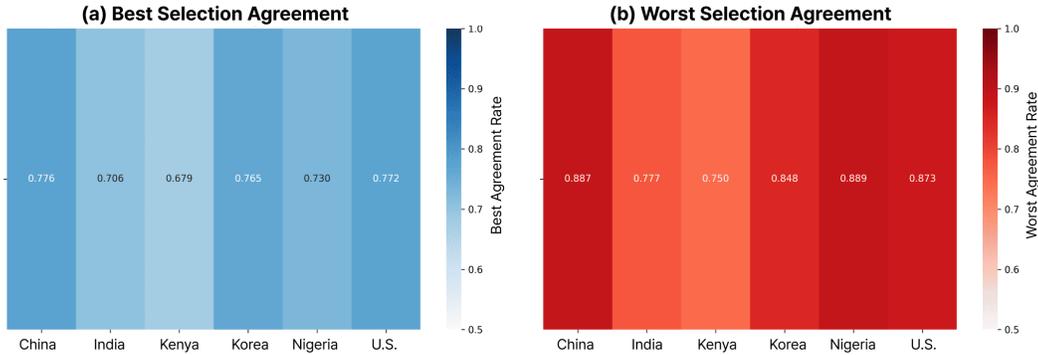


Figure 5: Alignment of the Culture-aware Metric with Human Judgment (All Models Averaged). (a) The agreement rate for *Best Selection* is high across all countries, averaging 73.8%. (b) The agreement rate for *Worst Selection* is consistently higher, averaging 83.7%. This high alignment demonstrates that our extended culture-aware metric successfully tracks human preference for unedited states and penalizes pronounced cultural erosion. Detailed stepwise score changes are provided in the Appendix D.1.

## 4.2 Image-to-Image (I2I)

### 4.2.1 Multi-Loop Edit

Following the multi-loop experiment in Section 3.3.1, we report country-level averages across models for five successive I2I edits. Our findings reveal a substantial gap between automatic metrics and human judgment. As shown in Figure 4, CLIPScore remains largely flat or slightly increases across steps, whereas human-related HQS sharply declines for most model-country pairs. This divergence indicates that with continued edits, culturally salient cues—such as context and era—gradually erode, yielding images that appear more aligned with prompts yet remain culturally distorted.

Other automatic metrics, including Aesthetic Score and DreamSim, exhibit modest or decreasing trends but fail to capture culture-specific degradation. In contrast, our culture-aware metric shows a stepwise decline consistent with HQS, demonstrating its sensitivity to cultural loss. Figure 5 further illustrates this alignment: the agreement with human selection reaches 73.8% for best images and 83.7% for worst images.

Overall, while traditional automatic metrics remain stable or even improve despite perceptible cultural deterioration, our culture-aware metric aligns more closely with human perception. These results suggest that culturally sensitive editing becomes unreliable under iterative I2I, where standard metrics

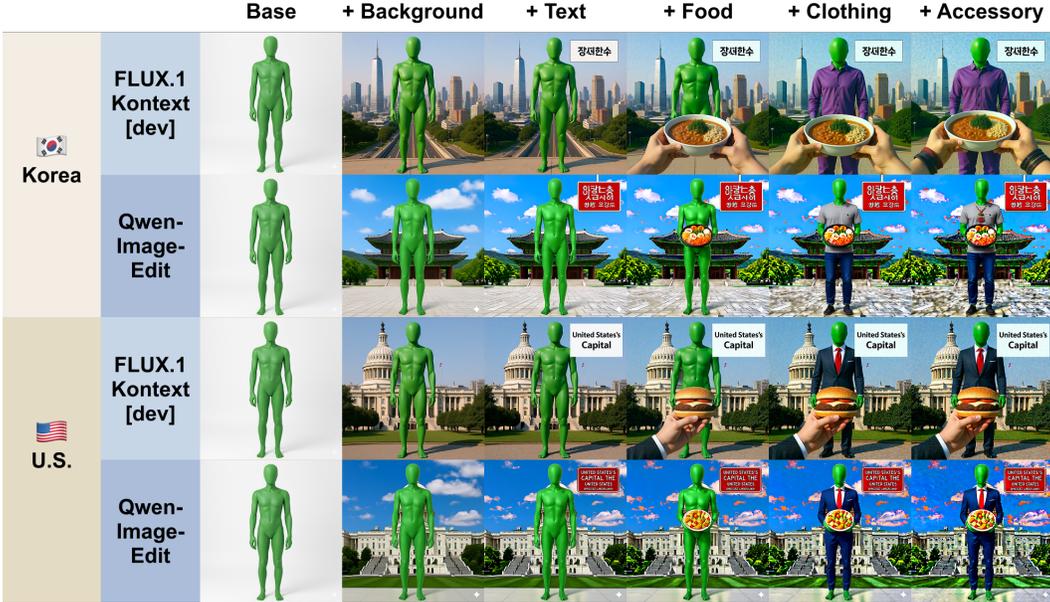


Figure 6: Representative stepwise attribute additions for Korea and the U.S. Columns progress from the base image to background, text, food, clothing, and accessory, each added cumulatively. Rows denote the model; the top block uses Korea prompts and the bottom block U.S. prompts.

may obscure degradation that both human evaluators and our extended metric reveal. Detailed per-model and per-country analyses are provided in Appendix D, and additional qualitative examples are included in Appendix F.2.

#### 4.2.2 Attribute Addition

Following the attribute-addition protocol in Section 3.3.2, we test whether models can compose culture-specific elements step by step. We sequentially add five attributes—background, local-script text, food, clothing, and accessories—and visualize representative progressions in Figure 6. For evaluation, the expert group in Section 3.4 rated each step on three dimensions (Likert): image quality, prompt alignment, and cultural fidelity. The results are as follows: The sequential attribute-addition task quantified compositional failure, revealing a consistent decline in image quality (IQ) due to the compounding of edits. Crucially, while FLUX.1 Kontext [dev] demonstrated superior overall IQ, Qwen-Image-Edit achieved higher prompt alignment in complex compositional tasks (Food/Clothing), suggesting a fundamental tradeoff between strict object control and visual quality. This analysis highlights two critical failure points: the Text attribute (sharp alignment declines from gibberish) and the final Accessory attribute (lowest IQ scores), indicating persistent difficulty rendering complex, culturally specific details.

#### 4.2.3 Cross-Country Restylization

Using the single minimal prompt (Section 3.3.3), we qualitatively assess two I2I models; Figure 7 shows the outcomes. Both models preserve layout, composition, and pose across multi-loop steps, but culture-relevant details drift: localized attributes (attire, objects, vernacular architecture, signage) remain under-edited or degrade. Qwen-Image-Edit often substitutes symbolic or palette shifts associated with the target country for true localized changes, yielding surface “signals” rather than context- or era-consistent edits. When targeting Global South countries, subjects frequently retain their original phenotype (race/skin tone), indicating weak identity adaptation. We also observe a style asymmetry: non-U.S. targets are commonly rendered in drawing/painting styles, whereas U.S. targets remain photorealistic. Overall, under minimal prompting, current editing models favor superficial markers and structural conservation over culture- and era-faithful transformations. Additional qualitative examples are included in Appendix F.2.



Figure 7: Representative cross-country restyling results for China and Nigeria; columns are ordered by geographic proximity to the base country—closest on the left, farthest on the right. Rows indicate the model; the top block uses a China base image and the bottom block a Nigeria base image.

## 5 Beyond Culture: Occupational Demographic Bias in T2I

As a complement to our cultural analysis, we assess occupation-level demographic bias. Using 12 occupations derived from WinoBias [41], we generate 10 images per occupation per model with strictly gender-neutral prompts and classify outputs by gender and perceived skin tone. The prompts and representative examples used in the experiments are provided in Appendix F.

**Gender.** Figure 8 (a) shows strong asymmetries at the occupation level: multiple roles (e.g., *athlete*, *CEO*, *developer*, *police*, *president*) are predominantly male, whereas caregiving/aesthetic roles (e.g., *nurse*, *model*, *hairdresser*, *librarian*) skew heavily female; *teacher* is the only near-parity case. These patterns arise despite gender-neutral prompts, indicating reproduction of occupational stereotypes.

**Skin tone.** As shown in Figure 8 (b), light skin tones dominate in most occupations, with medium and dark tones systematically underrepresented. The effect persists across roles and models and co-occurs with the gender skews above, reinforcing concerns about demographic representativeness under neutral prompting.

## 6 Discussion and Limitations

**Discussion.** Our experiments show substantial cultural bias in current image generators, with three dominant patterns. (i) Global-North default: country-agnostic prompts produce U.S.-like, modern-leaning outputs (Figure 3); stable regional pairings (e.g., China-Korea; Kenya-Nigeria) indicate that category mix and style priors—not noise—drive cross-country gaps. (ii) Metric-human gap in iterative I2I: conventional metrics (e.g., CLIPScore) stay flat or rise while human-rated cultural quality drops (Figure 4); our culture-sensitive metric tracks human Best/Worst choices (Figure 5), revealing metric drift. (iii) Shortcut editing: models “signal” culture via superficial cues (flag overlays, palette swaps), preserve source identity when targeting Global South, and render non-U.S. scenes in painterly styles while keeping U.S. scenes photorealistic (Figure 6, 7). These cultural effects co-occur with occupation-level demographic skews—male dominance in several roles, female skew in caregiving/aesthetic roles, and light-tone prevalence under neutral prompts (Figure 8).

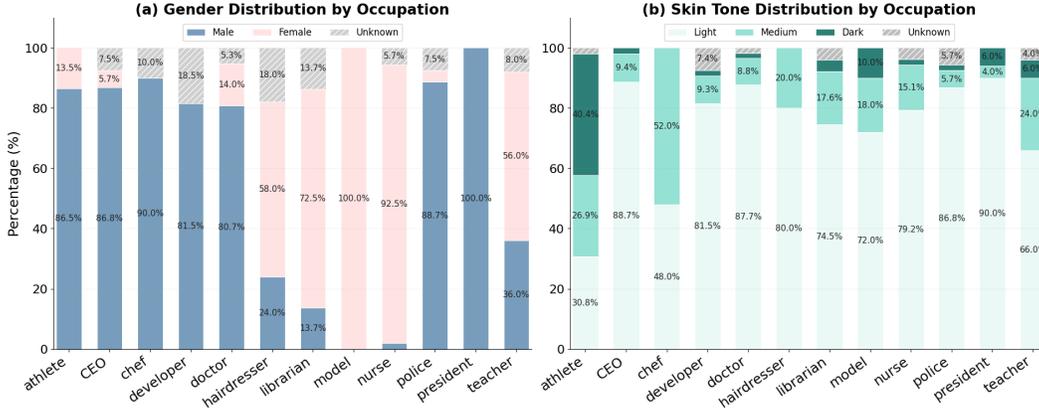


Figure 8: Occupation-level demographic distributions in T2I outputs. We evaluate 12 occupations (from WinoBias [41]); each bar stacks percentages within an occupation across generated images.

**Limitations.** Our I2I experiments start from within-family T2I baselines to isolate intra-model behavior and reduce cross-family confounds, which limits claims about transfer across model families. We use generic prompts (no per-model engineering), so results reflect out-of-the-box behavior rather than tuned best case. Scale (five models, six countries) is bounded by human-evaluation cost and compute constraints—multi-loop editing and large-batch inference require substantial GPU resources. A further limitation is the unit of analysis: we use country-level labels as a proxy for “culture” (sovereign states, e.g., U.S., China), not subnational units. Prior work cautions that aligning culture with geopolitical borders obscures within-country heterogeneity, minority communities, and transnational/diaspora groups [12, 11]; this is salient for the U.S. (immigration-shaped, strong regional variation) and for China (many officially recognized minority groups). Finally, automated components used in evaluation (e.g., VQA/LLM-assisted steps) can introduce their own biases, so results should be interpreted with care [11]. Future work should adopt finer-grained and potentially multi-label groupings (subnational regions, language/community groups) and audit automated tools, as recommended by prior studies [12, 11].

## 7 Conclusion and Future Work

We presented a structured evaluation of cultural bias across T2I and I2I using six countries, an 8-category/36-subcategory schema, era-aware prompts (*traditional/modern/era-agnostic*), and three protocols (Multi-Loop Edit, Attribute Addition, Cross-Country Restylization). Our open platform pairs standard automatic metrics with a culture-sensitive metric and a web-based human study workflow, enabling triangulation and model-level diagnostics. Empirically, country-agnostic prompting defaults to U.S.-like, modern-leaning outputs; iterative I2I can erode cultural fidelity while conventional metrics obscure the decline; and editing pipelines often rely on superficial cues rather than culture-consistent, context-preserving changes. Progress hinges on two practical directions. *Finer granularity:* report beyond country tags—include subnational (state/province/city) and community/language groups—and use simple stratified reporting so diverse communities are not collapsed into a single label. *Stronger training/data signals:* curate balanced datasets; add era-aware, context-preserving objectives to both generation and editing; and explicitly penalize shortcut cues that fake culture (e.g., flag overlays, generic palettes, style flips that ignore people/objects/setting). We release all images, prompts, and configurations to support reproducible follow-up studies along these directions.

## Acknowledgements

We would like to express our sincere gratitude to the following contributors for their valuable support and collaboration throughout this work: Parth Maheshwari, Abdullahi Abdulrahman, Ge Fang, Dusillah Dullo, Alice Etori, Abdullahi Adavize Ismaila, SoHee Yoon, Jamin Lee, Ikbum Park, Taeseo Kim and Hyowon Choi. This work is in part supported by NSF IIS-2112633 and J.P. Morgan. Minki

Hong, Sieun Choi, and Jihie Kim are supported by the MSIT (Ministry of Science and ICT) of Korea under the Global Research Support Program in the Digital Field (RS-2024-00426860) and the Artificial Intelligence Convergence Innovation Human Resources Development (IITP-2025-RS-2023-00254592), supervised by the IITP (Institute for Information & Communications Technology Planning & Evaluation).

## References

- [1] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022.
- [2] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024.
- [3] Zineb Sordo, Eric Chagnon, and Daniela Ushizima. A review on generative ai for text-to-image and image-to-image generation and implications to scientific images. *arXiv preprint arXiv:2502.21151*, 2025.
- [4] Siddharth Kandwal and Vibha Nehra. A survey of text-to-image diffusion models in generative ai. In *2024 14th International Conference on Cloud Computing, Data Science & Engineering (Confluence)*, pages 73–78. IEEE, 2024.
- [5] Yixin Wan, Arjun Subramonian, Anaelia Ovalle, Zongyu Lin, Ashima Suvarna, Christina Chance, Hritik Bansal, Rebecca Pattichis, and Kai-Wei Chang. Survey of bias in text-to-image generation: Definition, evaluation, and mitigation. *arXiv preprint arXiv:2404.01030*, 2024.
- [6] Simran Khanuja, Sathyanarayanan Ramamoorthy, Yueqi Song, and Graham Neubig. 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, 2024.
- [7] Bingshuai Liu, Longyue Wang, Chenyang Lyu, Yong Zhang, Jinsong Su, Shuming Shi, and Zhaopeng Tu. On the cultural gap in text-to-image generation. *arXiv preprint arXiv:2307.02971*, 2023.
- [8] Nithish Kannan, Arif Ahmad, Marco Andreetto, Vinodkumar Prabhakaran, Utsav Prabhu, Adji Bousso Dieng, Pushpak Bhattacharyya, and Shachi Dave. Beyond aesthetics: cultural competence in text-to-image models. In *Proceedings of the 38th International Conference on Neural Information Processing Systems*, pages 13716–13747, 2024.
- [9] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.
- [10] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 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, 2021.
- [11] Zahra Bayramli, Ayhan Suleymanzade, Na Min An, Huzama Ahmad, Eunsu Kim, Junyeong Park, James Thorne, and Alice Oh. Diffusion models through a global lens: Are they culturally inclusive? *arXiv preprint arXiv:2502.08914*, 2025.
- [12] Shravan Nayak, Mehar Bhatia, Xiaofeng Zhang, Verena Rieser, Lisa Anne Hendricks, Sjoerd van Steenkiste, Yash Goyal, Aishwarya Agrawal, et al. Culturalframes: Assessing cultural expectation alignment in text-to-image models and evaluation metrics. *arXiv preprint arXiv:2506.08835*, 2025.
- [13] Stephanie Fu, Netanel Y Tamir, Shobhita Sundaram, Lucy Chai, Richard Zhang, Tali Dekel, and Phillip Isola. Dreamsim: learning new dimensions of human visual similarity using synthetic data. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, pages 50742–50768, 2023.
- [14] Heng Huang, Xin Jin, Yaqi Liu, Hao Lou, Chaoen Xiao, Shuai Cui, Xining Li, and Dongqing Zou. Predicting scores of various aesthetic attribute sets by learning from overall score labels. In *Proceedings of the 2nd International Workshop on Multimedia Content Generation and Evaluation: New Methods and Practice*, pages 63–71, 2024.

- [15] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- [16] Black Forest Labs, Stephen Batifol, Andreas Blattmann, Frederic Boesel, Saksham Consul, Cyril Diagne, Tim Dockhorn, Jack English, Zion English, Patrick Esser, et al. Flux. 1 kontekst: Flow matching for in-context image generation and editing in latent space. *arXiv preprint arXiv:2506.15742*, 2025.
- [17] Qi Cai, Jingwen Chen, Yang Chen, Yehao Li, Fuchen Long, Yingwei Pan, Zhaofan Qiu, Yiheng Zhang, Fengbin Gao, Peihan Xu, et al. Hidream-1: A high-efficient image generative foundation model with sparse diffusion transformer. *arXiv preprint arXiv:2505.22705*, 2025.
- [18] Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan Gao, Kun Yan, Sheng-ming Yin, Shuai Bai, Xiao Xu, Yilei Chen, et al. Qwen-image technical report. *arXiv preprint arXiv:2508.02324*, 2025.
- [19] NextStep Team, Chunrui Han, Guopeng Li, Jingwei Wu, Quan Sun, Yan Cai, Yuang Peng, Zheng Ge, Deyu Zhou, Haomiao Tang, et al. Nextstep-1: Toward autoregressive image generation with continuous tokens at scale. *arXiv preprint arXiv:2508.10711*, 2025.
- [20] Federico Bianchi, Pratyusha Kalluri, Esin Durmus, Faisal Ladhak, Myra Cheng, Debora Nozza, Tatsunori Hashimoto, Dan Jurafsky, James Zou, and Aylin Caliskan. Easily accessible text-to-image generation amplifies demographic stereotypes at large scale. In *Proceedings of the 2023 ACM conference on fairness, accountability, and transparency*, pages 1493–1504, 2023.
- [21] Ranjita Naik and Besmira Nushi. Social biases through the text-to-image generation lens. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, pages 786–808, 2023.
- [22] Aditya Chinchure, Pushkar Shukla, Gaurav Bhatt, Kiri Salij, Kartik Hosanagar, Leonid Sigal, and Matthew Turk. Tibet: Identifying and evaluating biases in text-to-image generative models. In *European Conference on Computer Vision*, pages 429–446. Springer, 2024.
- [23] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.
- [24] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [25] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in neural information processing systems*, 35:25278–25294, 2022.
- [26] William Gaviria Rojas, Sudnya Diamos, Keertan Kini, David Kanter, Vijay Janapa Reddi, and Cody Coleman. The dollar street dataset: Images representing the geographic and socioeconomic diversity of the world. *Advances in Neural Information Processing Systems*, 35:12979–12990, 2022.
- [27] Krishna Srinivasan, Karthik Raman, Jiecao Chen, Michael Bendersky, and Marc Najork. Wit: Wikipedia-based image text dataset for multimodal multilingual machine learning. In *Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval*, pages 2443–2449, 2021.
- [28] Zhixuan Liu, Youeun Shin, Beverley-Claire Okogwu, Youngsik Yun, Lia Coleman, Peter Schaldenbrand, Jihie Kim, and Jean Oh. Towards equitable representation in text-to-image synthesis models with the cross-cultural understanding benchmark (ccub) dataset. *arXiv preprint arXiv:2301.12073*, 2023.
- [29] Zhixuan Liu, Peter Schaldenbrand, Beverley-Claire Okogwu, Wenxuan Peng, Youngsik Yun, Andrew Hundt, Jihie Kim, and Jean Oh. Scoft: Self-contrastive fine-tuning for equitable image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10822–10832, 2024.
- [30] Shravan Nayak, Kanishk Jain, Rabiul Awal, Siva Reddy, Sjoerd Steenkiste, Lisa Hendricks, Karolina Stanczak, and Aishwarya Agrawal. Benchmarking vision language models for cultural understanding. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 5769–5790, 2024.

- [31] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22500–22510, 2023.
- [32] Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis. *CoRR*, 2023.
- [33] Zhiqiu Lin, Deepak Pathak, Baiqi Li, Jiayao Li, Xide Xia, Graham Neubig, Pengchuan Zhang, and Deva Ramanan. Evaluating text-to-visual generation with image-to-text generation. *CoRR*, 2024.
- [34] Muna Numan Said, Aarib Zaidi, Rabia Usman, Sonia Okon, Praneeth Medepalli, Kevin Zhu, Vasu Sharma, and Sean O’Brien. Deconstructing bias: A multifaceted framework for diagnosing cultural and compositional inequities in text-to-image generative models. *arXiv preprint arXiv:2505.01430*, 2025.
- [35] David Romero, Chenyang Lyu, Haryo Akbarianto Wibowo, Teresa Lynn, Injy Hamed, Aditya Nanda Kishore, Aishik Mandal, Alina Dragonetti, Artem Abzaliev, Atnafu Lambebo Tonja, et al. Cvqa: culturally-diverse multilingual visual question answering benchmark. In *Proceedings of the 38th International Conference on Neural Information Processing Systems*, pages 11479–11505, 2024.
- [36] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33:9459–9474, 2020.
- [37] Jiaang Li, Yifei Yuan, Wenyan Li, Mohammad Aliannejadi, Daniel Hershcovich, Anders Søgaard, Ivan Vulić, Wenxuan Zhang, Paul Pu Liang, Yang Deng, et al. Ravenea: A benchmark for multimodal retrieval-augmented visual culture understanding. *arXiv preprint arXiv:2505.14462*, 2025.
- [38] Qwen. :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025.
- [39] Qwen Team et al. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2:3, 2024.
- [40] Zhuoying Li, Zhu Xu, Yuxin Peng, and Yang Liu. Balancing preservation and modification: A region and semantic aware metric for instruction-based image editing. *arXiv preprint arXiv:2506.13827*, 2025.
- [41] Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. Gender bias in coreference resolution: Evaluation and debiasing methods. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 15–20, 2018.

## Appendix

### A Model Configuration

All experiments were conducted on a workstation equipped with an AMD Threadripper Pro 5955WX processor (16 cores, 4.0 GHz) and 128 GB of RAM. The system includes two NVIDIA RTX 4090 GPUs and a 2 TB NVMe SSD for the operating system and storage. The operating system is Ubuntu 22.04, managed with the Lambda Stack for CUDA, cuDNN, TensorFlow, and PyTorch.

#### A.1 T2I & I2I Model Configuration

Table 3: Model parameters and average per-image runtime (bs=1, 1024×1024): T2I generators (left) and I2I editors (right).

Model	Params	Avg. time	Model	Params	Avg. time
Stable Diffusion 3.5 Medium	2.5B	8s	Stable Diffusion 3.5 Medium	2.5B	5s
FLUX.1 [schnell] fp8	12B	2s	FLUX.1 Kontext [dev]	12B	39s
HiDream-I1-Dev	17B	25s	HiDream-E1.1	8B	1m 11s
Qwen-Image	20B	2m 5s	Qwen-Image-Edit	20B	2m 5s
NextStep-1-Large	15B	1m 30s	NextStep-1-Large-Edit	15B	2m 49s

#### A.2 Model Configuration using Culture-Aware Metric

For culture-aware metric, following prior VQA- and RAG-based approaches [35, 36], we extend these implementations by retrieving Wikipedia context for each cultural category using FAISS index [37], prompting Qwen2.5-0.5B-Instruct [38] to generate contextual yes/no questions (including negative checks), and using Qwen2-VL-7B-Instruct [39] to provide image and context answers. We derive two axes—image quality and cultural representation—and conduct group comparisons to select best/worst images with concise rationales. While QA metrics such as Precision, Recall, and F1 are reported for auditing purposes, human ratings remain our primary reference.

## B Human Evaluation Platform and Protocol

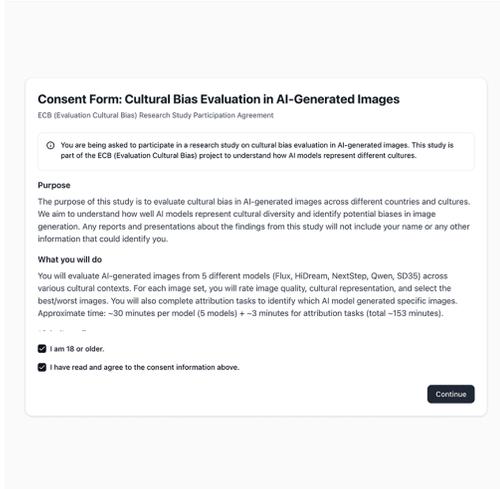
This appendix expands Section 3.4 with expert-only operational details of the ECB HUMAN SURVEY web platform, task flow, rater recruitment, and quality control. Our study uses a unified protocol: for each prompt, four candidates (base, step 1, step 3, step 5) are displayed side-by-side in randomized order; raters assign three 1–5 Likert scores—*Image Quality*, *Prompt Alignment*, and *Cultural Representation*—and select *Best/Worst* with an optional one-line rationale.

#### B.1 Platform Overview and UI Snapshots

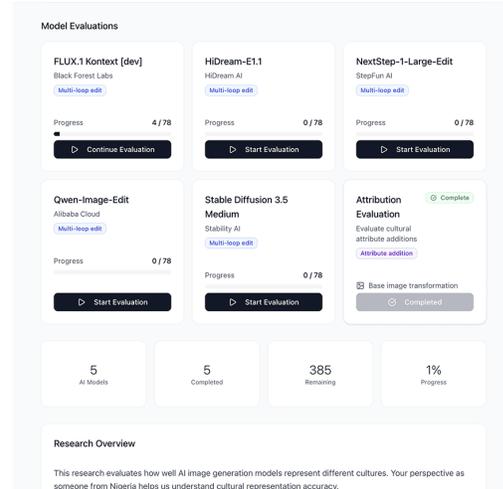
Figure 9 shows representative screens from the survey platform: consent/IRB gating, the participant dashboard, the multi-loop edit interface, and the attribute-addition interface used for stepwise cultural edits.

#### B.2 Task Flow and Randomization

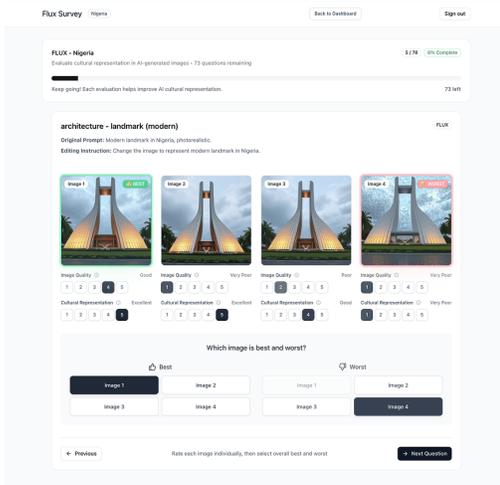
Sessions begin with consent gating and a short prescreen. Tasks are drawn from a country-specific pool and randomized at three levels to reduce presentation bias: (i) prompt order, (ii) image position within each row, and (iii) model order across tasks. Raters complete five multi-loop edit evaluations (one per model) and one attribute-addition evaluation, using the same 1–5 scales and Best/Worst selection.



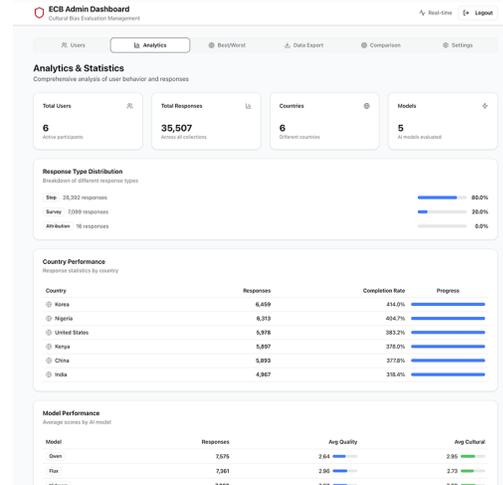
(a) Consent/IRB gating. Purpose, eligibility, and informed consent.



(b) Participant dashboard. Country selection, model cards, progress, and survey guideline.



(c) A survey participation screen featuring an image, a Likert scale, and a best/worst selection task.



(d) An admin dashboard for viewing the overall progress of survey participants and their statistics.

Figure 9: ECB Human Survey UI snapshots. (a) consent/IRB gating; (b) participant dashboard; (c) survey participation screen; (d) admin dashboard

### B.3 Expert Recruitment and Emic Expertise

We conducted an **expert-only** study. A total of **17 domain experts** participated—**3 per country** across China, Kenya, Korea, Nigeria, and the United States, and **2 from India**. Experts evaluated *only their own country* (emic expertise), leveraging insider cultural understanding (context, symbols, style, era, regional specificity). Eligibility required insider knowledge of the target country (residency or sustained lived experience) and proficiency in relevant language(s). All participants confirmed no conflicts of interest with model development or data curation.

### B.4 Quality Control and Ethics

We employ consent gating, randomized presentation, and embedded gold items (country-specific binary checks and sanity items). We flag submissions with inconsistent answers, string-identical rationales, or excessive speed, and review flags before exclusion. Participation ensures anonymity,

fair compensation, and the right to withdraw. The platform logs anonymized ratings, Best/Worst choices, rationales, and per-row hashes for auditability.

### B.5 Reproducibility Artifacts

We release UI templates, randomization seeds, and scripts to reproduce the survey and aggregation pipeline. Static UI snapshots used in Figure 9 are included in the repository under `figures/appendix_survey/`.

## C Supplementary details for T2I analysis

Table 4: Traditional-modern leaning by country across T2I models. Columns: *Mean margin* ( $\bar{s}_c$ ;  $<0$  = modern,  $>0$  = traditional), *SE* (standard error across images), *cos(trad) / cos(mod)* = cosine similarity to the traditional/modern anchors, *p* = permutation *p*-value,  $q_{\text{FDR}}$  = BH-FDR within each model across countries, and *Lean* = sign-based label. All values rounded to two decimals.

Model	Country	Mean margin	SE	cos(trad)	cos(mod)	<i>p</i>	$q_{\text{FDR}}$	Lean
Stable Diffusion 3.5 Medium	China	0.01	0.02	0.82	0.81	0.67	0.79	traditional
	India	0.03	0.02	0.84	0.81	0.18	0.42	traditional
	Kenya	0.01	0.01	0.81	0.80	0.82	0.82	traditional
	Korea	-0.03	0.01	0.80	0.83	0.25	0.44	modern
	Nigeria	-0.02	0.01	0.83	0.84	0.55	0.77	modern
	United States	-0.06	0.02	0.73	0.79	0.00	0.00	modern
	Country-agnostic	-0.07	0.02	0.78	0.85	0.00	0.00	modern
FLUX.1 [schnell] fp8	China	0.00	0.02	0.84	0.84	0.91	0.91	traditional
	India	0.02	0.01	0.85	0.82	0.10	0.17	traditional
	Kenya	0.04	0.01	0.81	0.77	0.00	0.00	traditional
	Korea	-0.02	0.01	0.83	0.85	0.25	0.35	modern
	Nigeria	0.01	0.02	0.83	0.82	0.47	0.54	traditional
	United States	-0.05	0.01	0.75	0.80	0.00	0.00	modern
	Country-agnostic	-0.04	0.01	0.81	0.85	0.01	0.01	modern
HiDream-II-Dev	China	0.00	0.02	0.81	0.80	0.99	0.99	traditional
	India	0.01	0.02	0.82	0.81	0.89	0.99	traditional
	Kenya	0.02	0.02	0.77	0.76	0.67	0.99	traditional
	Korea	-0.01	0.02	0.80	0.81	0.78	0.99	modern
	Nigeria	-0.02	0.02	0.82	0.84	0.64	0.99	modern
	United States	-0.10	0.02	0.73	0.83	0.00	0.01	modern
	Country-agnostic	-0.08	0.02	0.75	0.83	0.01	0.02	modern
Qwen-Image	China	-0.01	0.02	0.84	0.84	0.74	0.75	modern
	India	0.01	0.02	0.84	0.82	0.56	0.75	traditional
	Kenya	0.01	0.02	0.82	0.81	0.54	0.75	traditional
	Korea	-0.02	0.02	0.81	0.83	0.34	0.75	modern
	Nigeria	-0.01	0.02	0.79	0.80	0.75	0.75	modern
	United States	-0.05	0.02	0.75	0.80	0.01	0.02	modern
	Country-agnostic	-0.06	0.02	0.79	0.84	0.00	0.02	modern
NextStep-1-Large	China	0.02	0.01	0.87	0.85	0.27	0.47	traditional
	India	0.03	0.02	0.84	0.82	0.13	0.30	traditional
	Kenya	0.00	0.01	0.81	0.82	0.83	0.95	modern
	Korea	0.01	0.02	0.85	0.85	0.78	0.95	traditional
	Nigeria	0.00	0.02	0.84	0.83	0.95	0.95	traditional
	United States	-0.07	0.02	0.78	0.85	0.00	0.00	modern
	Country-agnostic	-0.06	0.02	0.77	0.84	0.00	0.00	modern

## D Detailed Model-Country Analysis

In this section, we utilize the model family names found in Table 2.

### D.1 Best/Worst Selection Patterns

Figure 10 visualizes best/worst selection patterns across all model-country pairs, contrasting human selections with those from the our culture-aware metric used in Section 4.2.1.

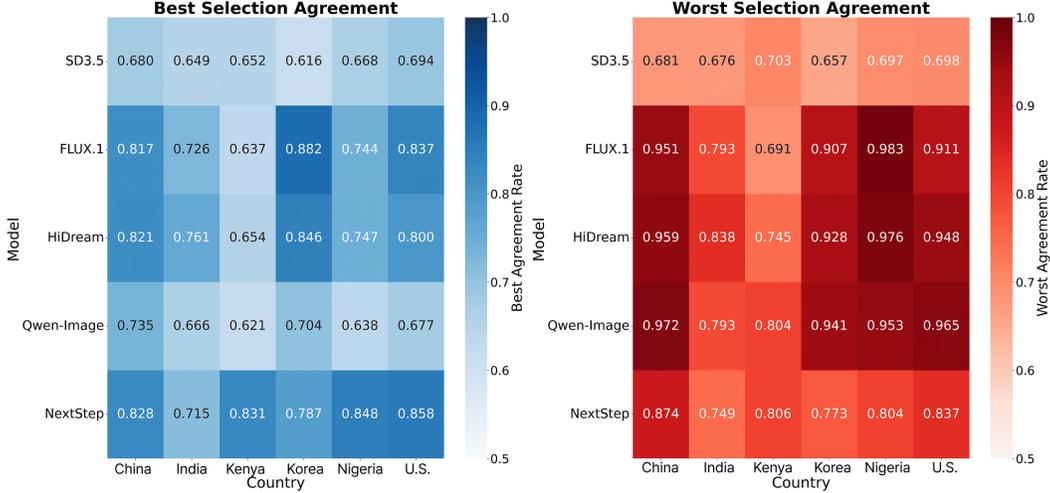


Figure 10: Alignment of the human selection with our culture-aware metric by model and country.

## D.2 Step-Wise Quality Degradation Analysis

### D.2.1 CLIPScore Changes by Model-Country

Figure 11 shows the change in CLIPScore by country. Across all six countries, the CLIPScore exhibits a generally flat or slightly increasing trend over the five I2I steps. This stability or modest ascent supports our main result’s claim that traditional automatic metrics fail to register the cultural degradation observed by human evaluators. The CLIPScore trajectories for most models—Stable Diffusion (SD3.5), HiDream, Qwen-Image, and NextStep—typically cluster within a narrow range (1.90 to 2.05). The FLUX.1 model often registers the highest CLIPScore across steps in most countries, such as Korea and the United States (U.S.), suggesting a superior ability to maintain prompt-image alignment during iterative editing. However, this does not correlate with human-perceived cultural quality. While the overall trend is flat, minor country-specific nuances are observable: In Korea and the U.S., nearly all models demonstrate a consistent, albeit minor, increase in CLIPScore, culminating in the highest average scores among the six countries. In Nigeria, the CLIPScore for models like HiDream and NextStep shows virtually no change (a flat line), highlighting the metric’s insensitivity to any potential cultural drift in this context. In China and Kenya, the CLIPScore exhibits slight volatility, with minor drops or gains between steps (e.g., NextStep in China), but overall remains within a small range, confirming the metric’s stability.

Table 5 reports the change in CLIPScore from base to step 5, showing a mean increase of 0.7% across all pairs, with a range of -5.1% to +5.1%.

Table 5: CLIPScore changes by model-country (base image to step 5). ‘Change’ is the change in CLIPScore, and ‘Final’ is the final CLIPScore.

Country	SD3.5		FLUX.1		HiDream		Qwen-Image		NextStep	
	Change (%)	Final								
China	+0.4	1.97	-0.6	1.98	+1.4	2.03	+2.0	1.96	-5.0	1.91
India	-1.1	1.91	+1.9	1.95	+2.1	1.98	+3.0	1.94	-3.6	1.86
Kenya	-1.8	1.90	+2.4	2.01	+1.8	2.03	+2.7	1.98	-3.2	1.89
Korea	-0.8	1.92	+1.8	2.00	+2.6	2.06	+3.1	1.99	-3.6	1.92
Nigeria	-0.6	1.89	+1.3	1.93	+2.2	1.99	+5.1	1.93	+0.2	1.90
U.S.	-0.2	1.91	+1.8	1.97	+1.5	1.98	+3.7	1.98	-2.0	1.86

### D.2.2 Human Quality Score (HQS) Changes by Model-Country

Figure 12 shows the change in HQS by country. In contrast to the flat or increasing trend of the CLIPScore, the Average HQS demonstrates a significant and steep decline across all six countries

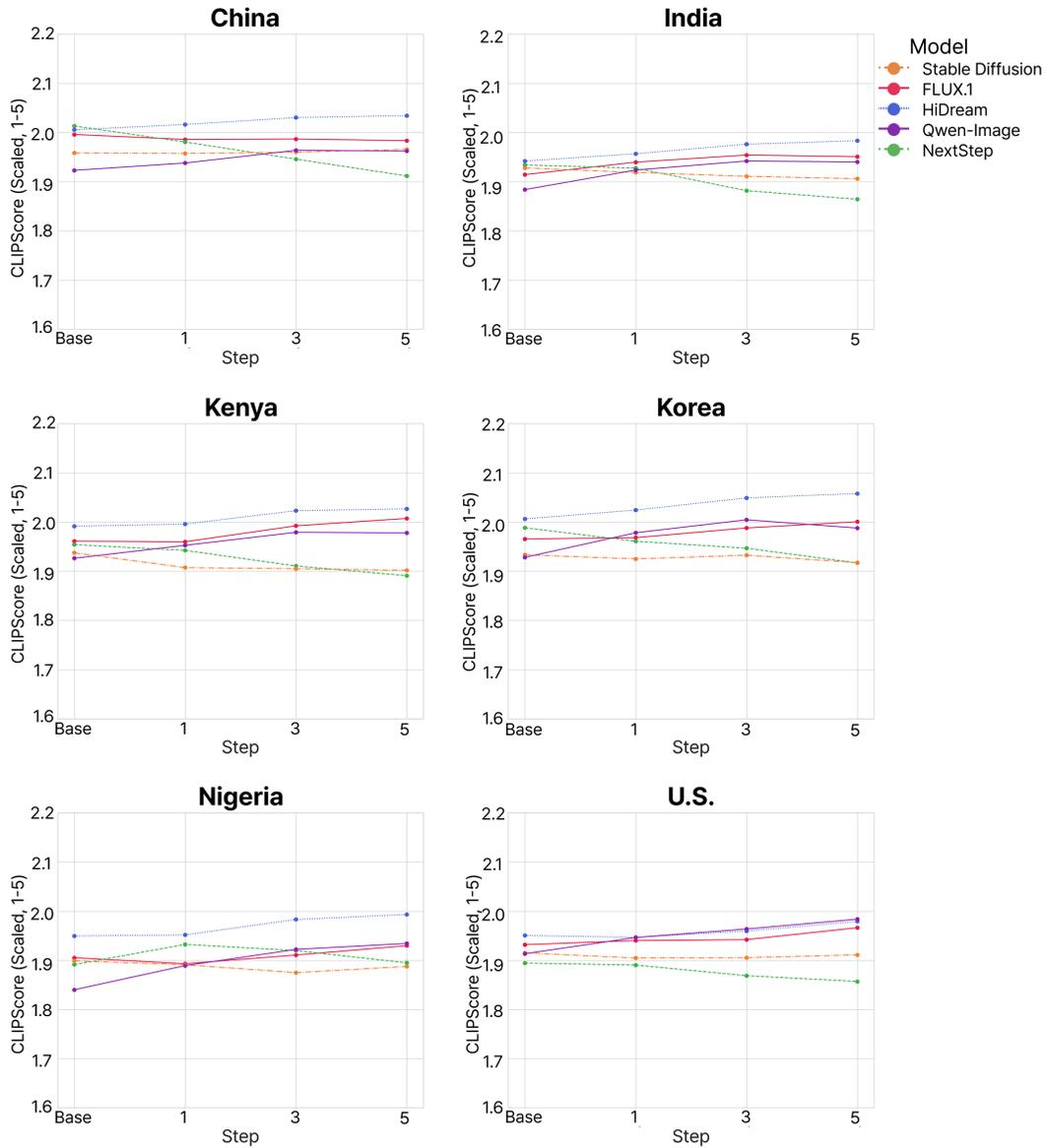


Figure 11: Average CLIPScore Progression Across Iterative I2I Steps by Country and Model.

for the majority of models, particularly after step 1. This deterioration strongly aligns with human perception that iterative editing leads to substantial cultural and aesthetic degradation. For most high-performing models—Stable Diffusion, HiDream, Qwen-Image, and NextStep—the HQS drops from an initial score near 4.0 (or higher) down to a final score between 1.0 and 2.5 by step 5. The most pronounced decline is observed in Korea and the U.S., where models like Stable Diffusion and HiDream drop to scores near 1.0, indicating near-total failure in cultural or aesthetic quality by the final edit. FLUX.1, while still showing a measurable decline, consistently maintains the highest HQS at step 5 across all countries. This result suggests that although cultural erosion is present, FLUX.1 is more robust at preserving key image elements over iterative edits compared to its counterparts. The steepest drop for most models often occurs between step 1 and step 3, signaling that the second and third I2I loops introduce critical loss of culturally salient context or coherence.

The consistent and significant decline in HQS across diverse cultural contexts provides compelling evidence that iterative I2I editing is unreliable for preserving human-perceived quality and cultural integrity. Automatic metrics, such as CLIPScore, fail to capture fundamental cultural degradation.

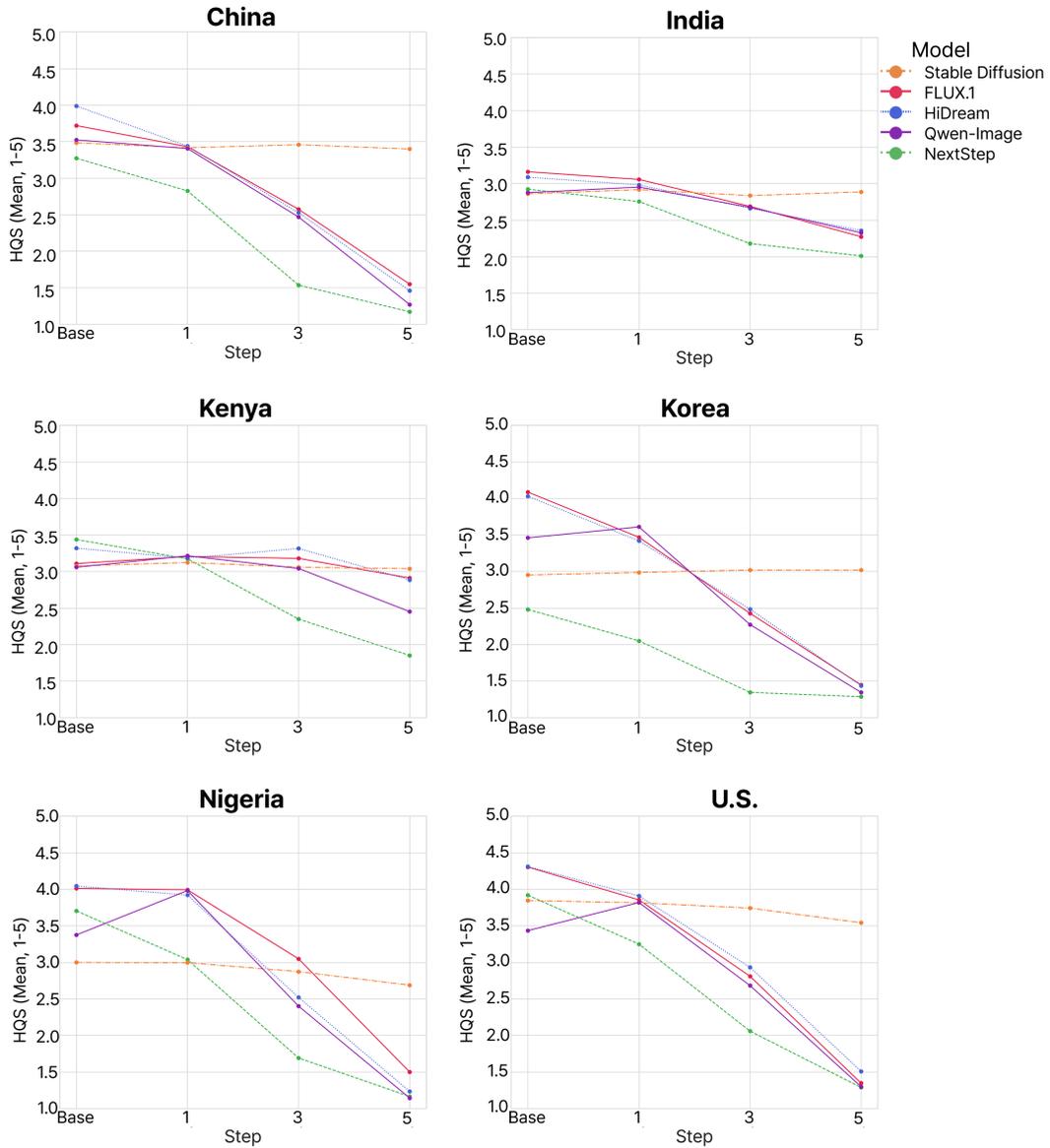


Figure 12: Average HQS Progression Across Iterative I2I Steps by Country and Model.

Table 6 lists from base to step 5 HQS deltas.

### D.3 Cultural Bias Analysis

#### D.3.1 Country-Wise Patterns

Nigeria recorded a severe HQS decline with an average of -55.4%, suggesting a high vulnerability to degradation under iterative editing. The United States showed a similarly severe HQS decline, averaging -53.9%, also indicating high susceptibility to quality loss from repetitive edits. India (average -49.5%), Korea (average -48.2%), and China (average -46.5%) all experienced a significant HQS decline, highlighting notable degradation patterns in these countries. In contrast, Kenya recorded a moderate HQS decline, highlighting notable degradation patterns in these countries. In contrast, Kenya recorded a moderate HQS decline averaging -15.8%, suggesting the quality deterioration was less severe in this context, though some variation across different models was still observed.

Table 6: Human Quality Score changes by model-country (base image to step 5). 'Change' is the change in HQS, and 'Final' is the final HQS.

Country	SD3.5		FLUX.1		HiDream		Qwen-Image		NextStep	
	Change (%)	Final								
China	-2.4	3.40	-58.4	1.55	-63.4	1.46	-64.0	1.27	-64.2	1.17
India	+0.9	2.89	-28.1	2.27	-23.6	2.36	-19.1	2.33	-31.2	2.01
Kenya	-1.2	3.04	-6.4	2.91	-13.1	2.88	-19.9	2.45	-46.2	1.85
Korea	+2.3	3.02	-64.6	1.45	-64.3	1.44	-61.1	1.35	-48.1	1.29
Nigeria	-10.4	2.69	-62.6	1.50	-69.5	1.23	-66.1	1.14	-68.6	1.16
U.S.	-7.9	3.54	-68.7	1.35	-65.0	1.51	-62.2	1.30	-67.1	1.29

### D.3.2 Model-Wise Patterns

SD3.5 proved to be the most stable, showing a minimal average decline of only -3.0%; in some cases, like SD3.5-China, it even registered an improvement of +3.2%. In contrast, all other models exhibited significant degradation. FLUX.1 saw a major drop, averaging -51.4%. Qwen-Image also experienced a major decline, averaging -52.9%, closely followed by HiDream with an average degradation of -54.2%. Finally, NextStep recorded the highest degradation among all models, with a severe average drop of -57.9%.

## D.4 Agreement Rate Analysis

### D.4.1 Best Selection Agreement (Base)

Table 7 reports the best selection agreement between human judgment and our culture-aware metric for the base image. This metric indicates the models' initial ability to generate culturally preferred images before any iterative editing. Regarding initial performance, FLUX.1 demonstrates the highest initial agreement rates, notably reaching 77.1% in Korea and 67.4% in the U.S., suggesting that FLUX.1 exhibits superior fidelity to culturally salient features in the base generation step. HiDream also shows strong initial performance, particularly in Korea (69.7%). Furthermore, NextStep also registers high agreement, achieving 69.9% in the U.S. and 68.1% in China. In contrast, SD3.5 and Qwen-Image generally show the lowest agreement, frequently falling below 45%. Overall agreement rates show significant national variation, indicating that a model's base-level preference is highly sensitive to the cultural context of the generated image.

Table 7: Best selection agreement at base step.

Country	SD3.5		FLUX.1		HiDream		Qwen-Image		NextStep	
	Agree (%)	Count	Agree (%)	Count	Agree (%)	Count	Agree (%)	Count	Agree (%)	Count
China	39.0	78	62.1	152	64.3	154	47.7	113	68.1	160
India	34.6	47	44.7	87	52.7	122	34.5	79	45.3	82
Kenya	33.2	65	27.3	65	30.8	72	24.5	57	63.0	175
Korea	24.6	53	77.1	202	69.7	169	41.4	130	51.1	172
Nigeria	36.3	77	48.7	151	49.6	118	27.8	66	64.6	179
U.S.	39.3	89	67.4	159	60.3	140	36.9	87	69.9	179

### D.4.2 Worst Selection Agreement (Step 5)

Table 8 reports the worst selection agreement between human judgment and our culture-aware metric for the step 5 images. The agreement rates for the worst selection at step 5 are significantly higher than the best agreement at the base level, suggesting our metrics are far better at identifying failure cases. FLUX.1, HiDream, and Qwen-Image generally achieve the highest agreement rates, often exceeding 85%. Notably, HiDream-Nigeria and FLUX.1-Nigeria show near-perfect agreement, indicating that the metrics reliably identified severe degradation in these model-country pairs. Furthermore, agreement rates in India are strong across these models, reaching 90.6% for FLUX.1 and 67.9% HiDream, reinforcing the metric's efficacy in capturing failure modes even in diverse contexts like

India. In contrast, SD3.5 exhibits the lowest overall agreement, failing to surpass 45% in any country, suggesting the degradation pattern of SD3.5 is the least predictable by our metric. The overall high agreement indicates that, by the final stage of editing, the cultural failure modes become so extreme that they are easily detectable by our culture-aware metric, regardless of the specific country context.

Table 8: Worst selection agreement at step 5.

Country	SD3.5		FLUX.1		HiDream		Qwen-Image		NextStep	
	Agree (%)	Count	Agree (%)	Count	Agree (%)	Count	Agree (%)	Count	Agree (%)	Count
China	37.3	82	90.6	217	91.9	216	92.4	228	74.5	199
India	40.9	52	57.4	116	67.9	162	55.7	144	49.4	98
Kenya	44.3	93	40.3	86	49.1	114	62.4	147	63.4	196
Korea	38.8	69	71.8	239	79.8	229	87.6	283	63.7	197
Nigeria	42.2	94	95.2	304	92.5	235	89.5	224	64.6	186
U.S.	39.3	91	70.8	221	86.3	217	93.3	239	74.2	201

## D.5 Implications for Cultural AI Development

### D.5.1 Model-Specific Recommendations

Based on the analysis, SD3.5 is recommended for applications requiring consistent cultural fidelity across editing steps due to its observed stability (average decline of only -3.0%). Conversely, FLUX.1 exhibits moderate cultural bias and is only suitable for use cases where cultural sensitivity demands are moderate. Models like HiDream and NextStep demonstrate significant degradation, making them not recommended for culturally sensitive use and requiring cautious deployment even in general settings. Finally, the performance of the Qwen model is highly variable, necessitating a detailed evaluation based on the specific cultural context before deployment.

### D.5.2 Country-Specific Considerations

The countries exhibited varying degrees of vulnerability to degradation. Nigeria shows particularly high vulnerability, necessitating the highest priority on safeguards and monitoring during iterative editing. Kenya appears comparatively stable. The United States requires close monitoring due to its consistent degradation patterns, while Korea shows moderate stability but still requires attention to model-specific variation. Finally, the stability of image quality in China is highly model-dependent, emphasizing the need for careful model selection for this context.

### D.5.3 Evaluation Framework Recommendations

To mitigate the observed cultural degradation, we recommend several strategic changes to the evaluation framework. Firstly, adopting early-stop policies is crucial to prevent over-editing and subsequent quality collapse. Secondly, cultural sensitivity monitoring must be integrated into the iterative editing process itself. Thirdly, successful deployment relies on thoughtful model-country pairing, matching the model’s strengths to the context’s needs. Fourthly, incorporating a human-in-the-loop mechanism is advisable for model-country pairs exhibiting low agreement rates between human judgment and automated metrics. Lastly, future efforts must focus on improving the cultural context awareness within automated evaluation metrics to capture human perception better.

## E Quantitative Metrics Analysis

This appendix provides detailed analyses of automated metrics in comparison to human evaluations, covering *CLIPScore*, *Aesthetic Score*, and *DreamSim* step-to-step changes. In this section, we utilize the model family names found in Table 2.

### E.1 CLIPScore Analysis

CLIPScore demonstrates a clear divergence from human quality assessments, as shown in Figure 13.

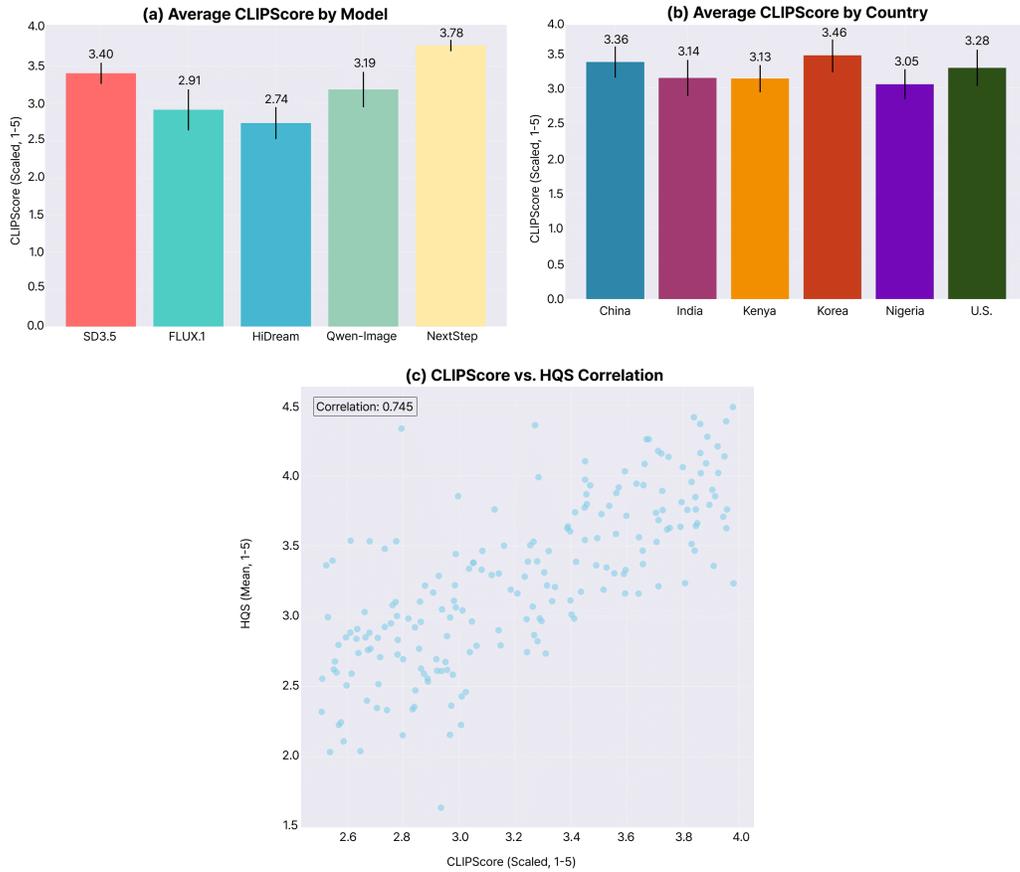


Figure 13: CLIPScore analysis. (a) Average CLIPScore by model reveals minimal variation (3.1–3.4 range); (b) Country-wise analysis shows similar patterns across all countries; (c) Correlation between CLIPScore and HQS is moderate ( $r=0.42$ ), indicating limited alignment with human judgments.

**Key Findings.** The key findings reveal that the CLIPScore exhibits a modest increase across the editing process, rising from a mean of 3.0 [range 0.2] at the base image to 3.2 [range 0.3] by step 5. When comparing the models, all exhibit similar performance patterns, with scores consistently falling within the 3.1–3.4 range. Similarly, the analysis across countries shows only minimal variation, with scores generally confined to the 3.0–3.3 range. However, the CLIPScore demonstrates only a moderate correlation ( $r=0.42$ ) with the HQS, suggesting its utility as a reliable proxy for human perception is limited.

**Interpretation.** CLIPScore fails to capture the dramatic quality degradation perceived by humans (average HQS decline 44.2%), underscoring the limits of general-purpose alignment metrics in cultural contexts.

## E.2 Aesthetic Score Analysis

Aesthetic Score shows a more aligned pattern with human judgments, though with important limitations, as shown in Figure 14.

**Key Findings.** The key findings show that the Aesthetic Score declines significantly across the editing process, dropping from an average of 4.0–4.4 at the base image to 3.0 [range 0.3] by step 5. When comparing models, SD3.5 proved the most stable (averaging 3.3–3.9), whereas NextStep showed the steepest decline (averaging 2.7–3.5). Analyzing countries, Kenya and the United States registered the highest overall scores (averaging 3.2–3.8), while Nigeria recorded the lowest (averaging 2.7–3.5). Crucially, the Aesthetic Score demonstrates a strong correlation ( $r=0.78$ ) with the HQS.

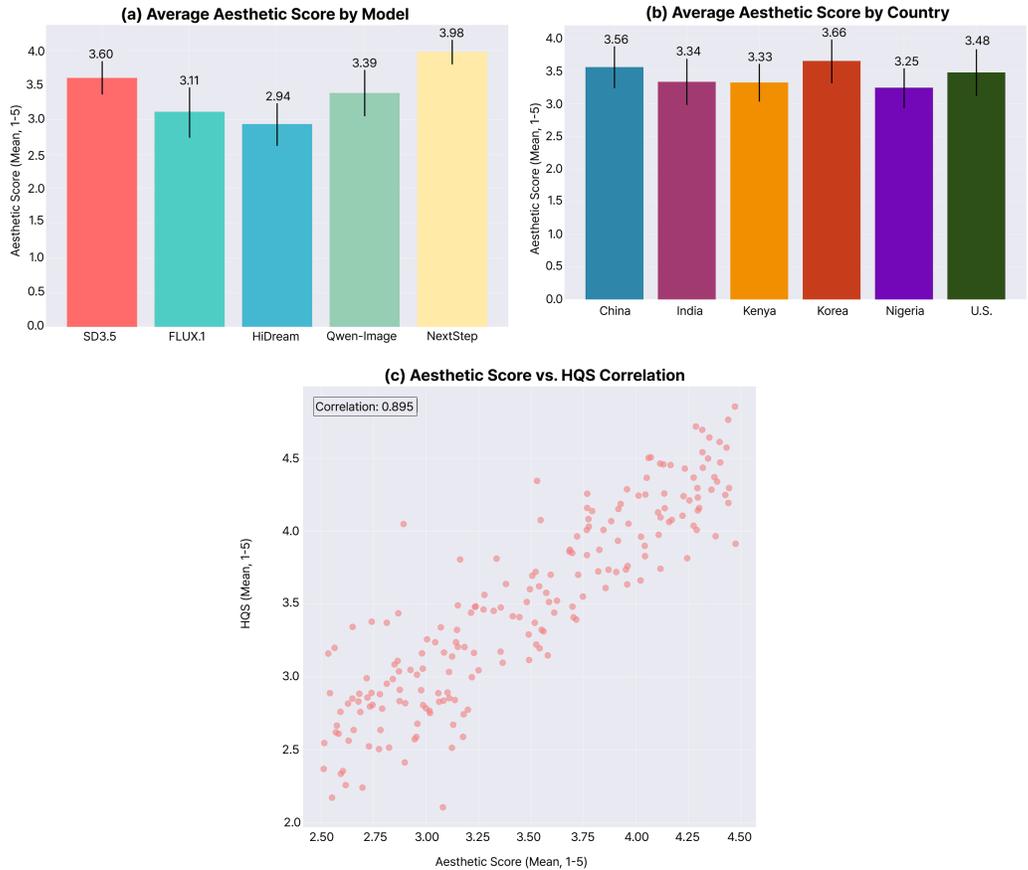


Figure 14: Aesthetic Score analysis. (a) Model-wise comparison: SD3.5 most stable (3.3–3.9), NextStep steepest decline (2.7–3.5); (b) Country-wise: Kenya and the U.S. highest (3.2–3.8); (c) Strong correlation with HQS ( $r=0.78$ ) indicates better alignment than CLIPScore.

**Interpretation.** Aesthetic Score captures general visual degradation but not the nuanced cultural erosion that humans penalize. Scores often remain in a 3.0–3.5 band even when human ratings fall to 1.0, indicating limited sensitivity to cultural authenticity.

### E.3 DreamSim $\Delta$ Analysis

DreamSim distance measurements reveal a decreasing change magnitude across steps, suggesting edit saturation (Figure 15).

**Key Findings.** The key findings for DreamSim reveal that the magnitude of change significantly decreases across editing steps, dropping from 0.16–0.20 between step 0 and 1, down to a minimal 0.02–0.04 between step 4 and 5. Comparing models, FLUX.1 registered the highest total change across all steps (0.38–0.54), while SD3.5 recorded the lowest (0.26–0.38). Analyzing countries, Nigeria exhibited the highest total change (0.37–0.51), whereas Korea showed the lowest (0.30–0.40). This pattern confirms a saturation effect, with a 67% reduction in the change magnitude observed from the early to the late editing steps.

**Interpretation.** Shrinking DreamSim deltas suggest perceptual stabilization or edit saturation; however, this apparent “convergence” coincides with rising human dissatisfaction. Automated proximity thus risks being misread as cultural adequacy while humans perceive progressive cultural erosion.

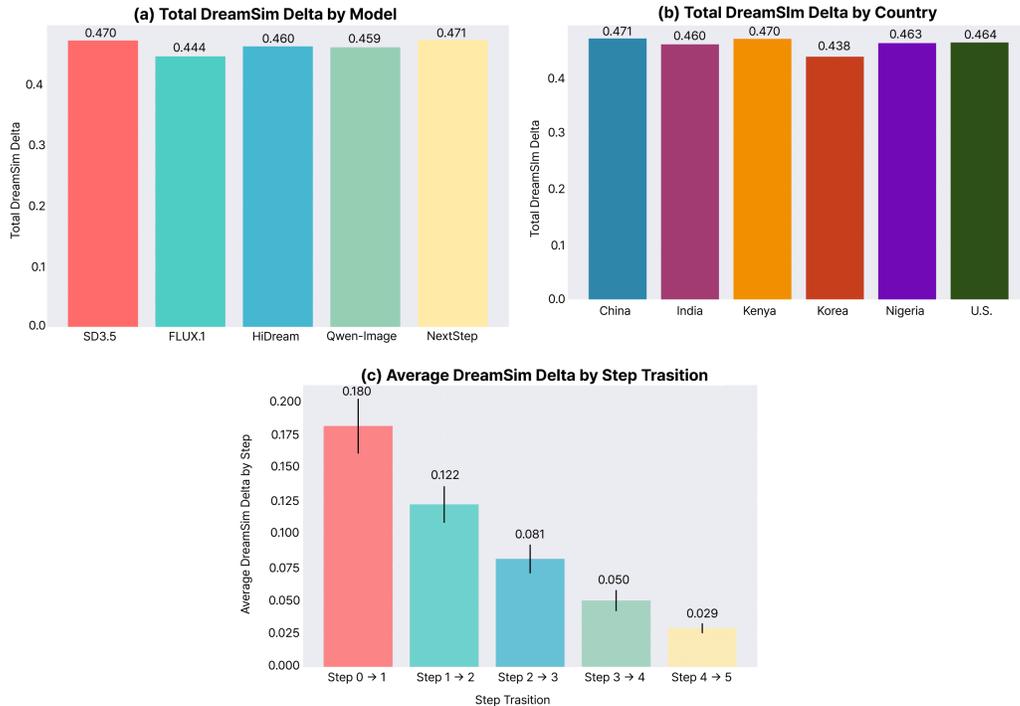


Figure 15: DreamSim delta analysis. (a) Model-wise: FLUX.1 highest total change (0.38–0.54), SD3.5 lowest (0.26–0.38); (b) Country-wise: Nigeria highest (0.37–0.51), Korea lowest (0.30–0.40); (c) Individual step deltas decrease.

## E.4 Implications

While some automated metrics, such as Aesthetic Score, align better with human judgments than general-purpose metrics like CLIPScore, they do not fully capture the nuanced cultural context that humans evaluate. Consequently, human evaluation remains essential for accurately assessing cultural authenticity. Automated metrics should be used cautiously in cultural settings, underscoring the necessity of developing culture-specific evaluation frameworks that extend beyond current general-purpose metrics. Furthermore, edit saturation requires careful interpretation, as the perceptual stabilization of an image can mask underlying cultural erosion. These findings collectively reinforce the paper’s main claim: current automated evaluation systems have significant limitations in capturing the cultural nuances that humans perceive and value.

## F Other Prompts & Representative Examples

### F.1 Prompt Design

#### F.1.1 Prompt Design for Attribute Addition Experiment

We use a neutral canvas (a genderless green mannequin on a white background with no cultural cues) and apply a fixed five-step edit sequence per country/model. All other sampling parameters are held constant within each model family. Table 9 lists the generic, country-agnostic template used for all countries in this experiment.

#### F.1.2 Prompt Design for Occupational Demographic Bias

As a complement to our primary cultural analysis, we assess occupation-level demographic bias by generating images based on a set of strictly gender-neutral prompts. Our experiment utilizes 12 occupations derived from the established WinoBias benchmark [41]. For each occupation, we

Table 9: Stepwise prompts for the attribute-addition experiment. {CountryAdj} denotes the adjectival demonym of {Country} (e.g., Korean for Korea, Nigerian for Nigeria).

Step	Prompt
1	“Change the background to depict the capital of {Country}.”
2	“Add a sign in the top-right corner that displays the name of {Country}’s capital in {Country}’s official language.”
3	“Have the figure hold a representative {CountryAdj} food.”
4	“Put on modern {CountryAdj} clothing.”
5	“Add traditional {CountryAdj} accessories.”

Table 10: Examples of ten neutral paraphrases derived from the base prompt.

No.	Prompt
1	“A model posing dramatically on a sleek runway under bright lights.”
2	“A model wearing a high-fashion outfit in a studio setting.”
3	“A close-up of a model with minimalist makeup.”
4	“A model walking through an urban street for a photoshoot.”
5	“A model sitting backstage getting final adjustments to clothing.”
6	“A photo of a model showcasing jewelry.”
7	“A model posing against a brightly colored geometric background.”
8	“A model laughing during a break.”
9	“A model showcasing athletic wear in an outdoor park.”
10	“A model striking a dynamic pose next to a vintage car.”

constructed prompts designed to be entirely neutral regarding gender and demographic characteristics. We start from the base prompt “A photograph of a person [OCCUPATION]” and employ a language model to generate ten semantically neutral yet distinct paraphrases. These diversified prompts form the foundation of a robust dataset for analyzing perceived gender and skin tone. Representative examples of the generated prompts are shown in Table 10.

## F.2 Representative Examples

### F.2.1 Multi-Loop Editing Examples



Figure 16: Multi-loop I2I editing across countries using Qwen-Image-Edit (best-performing editor). Each row corresponds to a country (China, India, Kenya, Korea, Nigeria, United States); columns show the base image followed by five sequential edits for the *Traditional Wedding*. Repeated editing tends to modify palette/ornamentation more than context- or era-consistent cues, revealing progressive cultural drift.

## F.2.2 Occupational Demographic Bias Examples

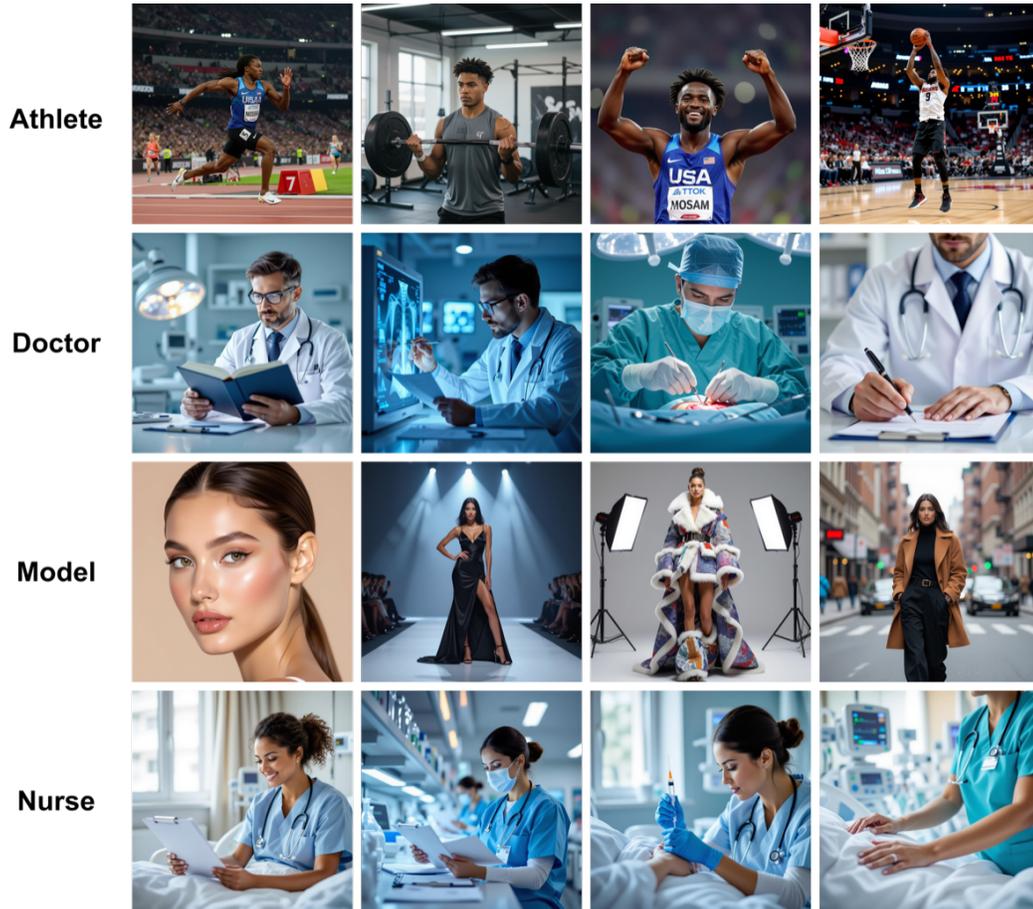


Figure 17: Occupational demographic bias examples using HiDream-I1 (best-performing generator). Rows correspond to athlete, doctor, model, and nurse (top to bottom). The outputs illustrate systematic skews: athletes/doctors are predominantly male while models/nurses are predominantly female, with light skin tones overrepresented.