# Deconstructing Bias: A Multifaceted Framework for Diagnosing Cultural and Compositional Inequities in Text-to-Image Generative Models

**Anonymous authors**
Paper under double-blind review

## Abstract

The transformative potential of text-to-image (T2I) models hinges on their ability to synthesize culturally diverse, photorealistic images from textual prompts. However, these models often perpetuate cultural biases embedded within their training data, leading to systemic misrepresentations. This paper benchmarks the Component Inclusion Score (CIS), a metric designed to evaluate the fidelity of image generation across cultural contexts. Through extensive analysis involving 2,400 images, we quantify biases in terms of compositional fragility and contextual misalignment, revealing significant performance gaps between Western and non-Western cultural prompts. Our findings underscore the impact of data imbalance, attention entropy, and embedding superposition on model fairness. By benchmarking models like Stable Diffusion with CIS, we provide insights into architectural and data-centric interventions for enhancing cultural inclusivity in AI-generated imagery. This work advances the field by offering a comprehensive tool for diagnosing and mitigating biases in T2I generation, advocating for more equitable AI systems.

## 1 Introduction

Synthetic image generation has emerged as a transformative computational paradigm, with diffusion models and GANs enabling photorealistic visual synthesis from textual or structured inputs. Systems like DALL·E 3 and Gemini exemplify this capability, driving revolutionary applications across creative industries, computer vision pipelines, and AI-assisted design. Built on transformative advancements in deep learning architectures demonstrate unprecedented capability in synthesizing photorealistic images from textual prompts. These models, built on transformer architectures (Rombach et al., 2022) and diffusion processes (Ho et al., 2020), are now integral to applications spanning creative industries, education, and cultural preservation.

However, as these models transition from research curiosities to production environments, fundamental challenges emerge: the same architectures that achieve unprecedented image fidelity systematically amplify societal biases encoded in their training corpora limiting global representational accuracy. This work introduces a rigorous evaluation framework utilizing the Components Inclusion Scores (CIS) to quantify understudied bias dimensions: (1) compositional fragility in multi-element synthesis and (2) contextual misalignment in culturally nuanced prompts

### 1.1 The Bias Amplification Challenge

Bias in text-to-image generation arises when models produce outputs that reflect and potentially amplify societal stereotypes present in their training data. These biases manifest in various forms, such as gender, skin tone, and cultural representations, leading to images that may not accurately or fairly depict the intended subjects. For instance, prompts describing "a traditional wedding" often

generate Western-style ceremonies, while non-Western cultural elements are underrepresented or misrepresented.

Our large-scale analysis of 2,400 generated images reveals three systemic failure modes in state-of-the-art T2I models:

1. **Compositional Fragility**: Models struggle to accurately combine marginalized cultural elements, leading to significant performance disparities compared to Western counterparts.

2. **Contextual Degradation**: The inclusion of historical and cultural contexts disproportionately reduces accuracy for non-Western concepts, indicating a bias in contextual fidelity.

3. **Order Sensitivity**: The sequence of elements within a prompt introduces performance instability, with significant variations in output quality depending on element ordering.

Detailed examples and quantitative results for these failure modes are presented in experiments section.

## 1.2 ROOT CAUSES OF BIAS

The systemic failures observed in T2I models stem from three interconnected technical limitations, each contributing to the amplification of cultural and compositional biases. These limitations are deeply rooted in the data, architecture, and optimization dynamics of modern generative models.

1. Training Data Imbalance: LAION-5B, the primary training corpus for many T2I models, contains 18× more Western cultural references than African/Asian artifacts (Schuhmann et al., 2022). This skew propagates through the generation pipeline, as shown by PCA analysis of latent embeddings (Fig. 2).

2. Architectural Limitations: Cross-attention layers exhibit 3.2× higher entropy for minority concept pairs ($H = 3.8$ vs. $H = 1.2$ for mainstream), correlating with omission/conflation errors ($r = -0.71$).

3. Embedding Superposition: Minority cultural concepts occupy overlapping latent dimensions (68% overlap vs. 22% for mainstream), a consequence of transformer models compressing rare tokens into shared parameter space (Elhage et al., 2021).

## 1.3 LIMITATIONS OF CURRENT APPROACHES

Existing bias mitigation strategies fail to address the multifaceted nature of cultural and compositional biases in T2I models. Below, we dissect these limitations across three dimensions, supported by empirical and theoretical evidence:

1. Surface-Level Interventions: Methods like dataset balancing (Li et al., 2022) and adversarial debiasing reduce overt stereotypes (e.g., "CEO" → male) but fail to address nuanced cultural misrepresentations. For instance, Bansal (2022) reduced gender bias in Stable Diffusion by 37% but reported no improvement in cultural accuracy for non-Western prompts.

2. Cultural Blindness: Studies like "Fair Diffusion" Friedrich et al. (2023a) focus on equalizing demographic attributes (e.g., skin tone distribution) but ignore contextual fidelity (e.g., traditional attire in cultural ceremonies).

3. Lack of Cross-Cultural Evaluation: Benchmarks such as BiasBench test only 5% of prompts on non-Western cultural concepts, leaving systemic underrepresentation unmeasured.

### 1.3.1 METRIC GAPS: THE PHANTOM OF OBJECTIVITY

Traditional evaluation metrics prioritize technical quality over fairness, creating a false sense of progress:

- Explicit vs. Implicit Bias: Current metrics like FairFace (Karkkainen & Joo, 2021) detect overt stereotypes (e.g., racial mis-classification) but miss implicit biases, such as the conflation of "Moroccan lanterns" with Chinese designs.

- Contextual Ignorance: CLIP-based metrics measure prompt-image alignment but fail to penalize cultural inaccuracies (e.g., a "Nigerian wedding" generated in a Gothic church)

- Static Evaluations: Benchmarks test single-concept prompts (e.g., "doctor"), ignoring compositional failures (e.g., "Indian scientist in a lab with traditional art").

### 1.3.2 Architectural Blind Spots: Symptomatic Solutions

State-of-the-art bias mitigation strategies often address surface symptoms rather than underlying architectural limitations. Prompt engineering, such as adding culturally specific terms ("traditional Ugandan design"), can improve CIS by 15% but requires manual intervention and fails to correct data imbalances, leading to inconsistent gains (±22% CIS variation) (Bianchi et al., 2023). Dataset filtering reduces overt stereotypes by 40% but unintentionally removes 68% of non-Western cultural references due to automated NSFW filters, causing a 52% CIS decline for marginalized prompts(Schuhmann et al., 2022). Adversarial training penalizes biased outputs but at the cost of model performance, increasing FID by 0.19 and reducing CIS by 0.33 (Zhang et al., 2024). Despite their effectiveness in mitigating immediate biases, these strategies do not fundamentally resolve the deeper architectural challenges that contribute to systemic inconsistencies in AI-generated content.

### 1.3.3 The Missed Nexus: Data, Architecture, and Culture

Current approaches overlook the interplay between data imbalance, transformer dynamics, and cultural semantics: Data-Centric Myopia: Methods like data augmentation add synthetic examples but ignore how minority embeddings are compressed via superposition. Architectural Rigidity: Post-hoc fixes (e.g., attention layer fine-tuning) fail to address cross-attention entropy spikes for minority pairs. Cultural Atomization: Treating cultural concepts as isolated tokens (e.g., "kimono") rather than contextual systems (e.g., "Japanese tea ceremony") leads to fragmented representations.

### 1.4 Our Contributions

We benchmark cultural bias in text-to-image generative models using the Component Inclusion Score (CIS), which integrates component inclusion, contextual alignment, and cultural fidelity to quantify disparities in generated outputs. Our analysis reveals significant performance gaps, with models underperforming on non-Western cultural prompts compared to Western-centric ones ($p < 0.001$). We identify underlying causes such as training data imbalances, elevated cross-attention entropy, and latent embedding superposition. These findings highlight critical shortcomings in current models and offer actionable insights for addressing biases through architectural and data-centric interventions, advancing fairness and inclusivity in generative models.

## 2 Related Work

### 2.1 Bias in Generative Models

Recent advances in text-to-image generative models, such as DALL·E 3, Stable Diffusion, and Gemini, have demonstrated remarkable image synthesis capabilities (Rombach et al., 2022; Ho et al., 2020). However, these systems often inherit and amplify societal biases present in their training data. For instance, LAION-5B—a primary training corpus for many of these models—has been shown to contain up to 18× more Western cultural references than African/Asian artifacts (Schuhmann et al., 2022), leading to cultural misrepresentations. Additional work in transformer dynamics (Elhage et al., 2021) and parameter-efficient fine-tuning (Houlsby et al., 2019) further highlights the challenges of aligning model architectures with culturally diverse representations.

### 2.2 Evaluation Metrics and Diagnostic Tools

Traditional evaluation metrics such as FID and CLIP-based scores primarily assess image quality and semantic alignment, often overlooking nuanced cultural and contextual inaccuracies. Efforts to mitigate bias have included surface-level interventions like dataset balancing (Li et al., 2022) and fairness-oriented datasets like FairFace (Karkkainen & Joo, 2021). However, these approaches

often miss deeper implicit biases such as contextual misalignment and compositional fragility. Additionally, recent studies such as "Fair Diffusion" by Friedrich et al. (2023b) have begun to address cultural representation issues, yet our CIS advances this line of research by providing a detailed quantification of both explicit and implicit biases in generated outputs.

### 2.3 ARCHITECTURAL DRIVERS OF BIAS

Our analysis reveals two intertwined mechanisms that amplify biases in generative models. First, *superposition* occurs when latent representations of rare tokens become overwritten by more dominant patterns, effectively compressing multiple cultural features into a shared embedding space. This phenomenon undermines the distinct representation of minority concepts, as detailed by (Elhage et al., 2021). Second, the observed non-monotonic error curve in our models aligns with the *double descent* phenomenon, where increasing model complexity can initially increase error rates before decreasing them. This effect particularly impacts the accurate representation of minority cultural elements due to phase transitions in training. Our findings indicate that the high overlap in embeddings for marginalized cultural concepts (68%)—compared to only 22% for mainstream concepts—directly correlates with a collapsed latent space. In such a space, diverse cultural elements are not distinctly represented, leading to conflation in generated imagery. Together, these architectural factors underscore the need for model design strategies that mitigate the adverse effects of data imbalance and latent embedding interference on representational fidelity.

## 3 METHODOLOGY

### 3.1 COMPONENT INCLUSION SCORE (CIS)

CIS is a quantitative metric designed to measure how accurately a generative model incorporates specified components from a prompt into the generated imageChen et al. (2023). In our study, CIS was used to evaluate biases in image generation when depicting subjects from both marginalized and non-marginalized countries, specifically in the categories of flags, monuments, vehicles, and food. Ideally, for each prompt containing key components—such as cultural artifacts, geographic references, or demographic attributes—the model should accurately render all these elements in the generated image. A higher CIS score indicates a model's ability to faithfully represent complex prompts without omitting critical components.

The CIS score for an individual image $I_{i,j}$ is calculated as:

$$S_{i,j} = \frac{L(\text{argmax}(\hat{p}_{i,j}))}{K},$$

where $L(\text{argmax}(\hat{p}_{i,j}))$ is the number of components successfully identified from the lookup table $L$ for the image $I_{i,j}$. The final CIS metric for a given number of components $K$ is computed as:

$$\text{CIS}_K = \frac{1}{M \cdot N} \sum_{i=1}^{M} \sum_{j=1}^{N} S_{i,j}.$$

The CIS metric serves as a robust indicator of how effectively the model retains and represents multiple elements from a prompt, allowing us to quantify any disparities in image generation for marginalized versus non-marginalized groups.

### 3.2 EXPERIMENTAL DESIGN

We classified prompts into four categories, each consisting of 100 distinct concepts, to evaluate the performance of text-to-image models:

1. Base Prompts: Single-concept prompts featuring well-known subjects (e.g., "Big Ben," "Taj Mahal").
2. Pair/Trio Prompts: Combinations of two or three distinct concepts (e.g., "Eiffel Tower + Vesak Lanterns," "Statue of Liberty + Diwali Lamps + Sombrero").

3. Contextual Prompts: Prompts with specific cultural or historical contexts (e.g., "Moroccan market with traditional textiles," "Japanese tea ceremony in a zen garden").

4. Adversarial Prompts: Perturbed prompts designed to test model robustness by introducing incongruent elements (e.g., "Ancient Egyptian pyramid in New York City," "Futuristic samurai in a medieval European castle").

Models Evaluated:

In this study, we evaluate the following text-to-image models:

- Stable Diffusion v2.1:A diffusion-based generative model for creating images from text descriptions, widely recognized for its ability to produce high-quality outputs.
- SG161222/Realistic Vision V1.4: A model fine-tuned for photorealistic image generation, built upon the SG161222 architecture to enhance visual realism.
- Dreamlike-Art/Dreamlike Photoreal 2.0: A model designed for generating detailed and lifelike images, with a focus on high-fidelity photorealistic rendering.

Each model has undergone pre-training on large-scale image-text datasets, with configurations and parameters used according to their respective specifications. For consistency and reproducibility, the temperature setting was fixed at 0 for all models during evaluation.

Validation Protocol:

1. Automated Scoring: We employ CLIP and Mask R-CNN for objective evaluation of generated images. CLIP assesses the overall semantic similarity between the prompt and the generated image. Mask R-CNN identifies specific objects and their spatial relationships within the image. These scores are combined to form a comprehensive automated evaluation metric.

2. Architectural Analysis: We analyze attention maps to understand which parts of the image the model focuses on for different cultural contexts. Principal Component Analysis (PCA) is performed on the embedding space to visualize how different cultural concepts are represented in the model's latent space.

## 4 EXPERIMENTS

### 4.1 PERFORMANCE DISPARITIES

To evaluate how well each model captures cultural elements and adapts to different prompts, we analyze their performance across four key dimensions: cultural representation, compositional accuracy, contextual consistency, and robustness in historical and modern settings. Results are summarized in Table 2 (see Appendix A). Stable Diffusion v2.1 exhibited broad cultural representation but struggled with fine-grained differentiation. Realistic Vision V1.4 excelled in contextual consistency and historical accuracy but underperformed in blending distinct cultural elements. Dreamlike Photoreal favored Western-centric elements and showed the lowest performance in cross-cultural pairings and historical settings.

| Metric | Photorealism (↑) | Fairness Sensitivity (↑) | Cultural Nuance (↑) |
|---|---|---|---|
| FID (Heusel et al., 2017) | 0.92 | 0.12 | 0.08 |
| CLIP-Score (Radford et al., 2021) | 0.85 | 0.31 | 0.24 |
| CIS (Ours) | **0.88** | **0.79** | **0.68** |

Table 1: Normalized metric performance on 200 culturally diverse prompts (higher is better).

The results highlight the limitations of conventional evaluation metrics, which tend to favor photorealism at the expense of fairness and cultural inclusivity. Our findings suggest that models optimized solely for FID or CLIP-Score may reinforce cultural biases by underrepresented marginalized aesthetics, whereas CIS provides a more holistic evaluation framework.

## 4.2 ARCHITECTURAL ANALYSIS

To analyze the variations in cross-attention entropy across transformer layers, we observe a noticeable peak at layer 6, as shown in Figure 1 (see Appendix B). Additionally, the framework for the CIS metric is illustrated in Figure 2 (see Appendix B).This underscores the need for targeted interventions at these architectural layers to mitigate bias.

## 5 ANALYSIS AND DISCUSSION

### 5.1 ROOT CAUSES OF SYSTEMIC BIASES

Our analysis revealed systemic biases in text-to-image models driven by two primary factors. First, the LAION-5B dataset, despite its 5.85 billion image-text pairs, is culturally imbalanced: only 12.7% of non-Western cultural artifacts appear in at least five instances, versus 89% for Western artifacts. This disparity arises from CLIP filtering bias—using similarity thresholds of 0.28 for English and 0.26 for other languages that disproportionately filter out non-Western content—and from a skewed source distribution, with 78% of English-language pairs coming from North American and European domains compared to just 6% from African or Asian sources.

Second, architectural limitations in transformer-based models contribute to bias. Superposition, where overlapping embedding subspaces encode multiple concepts, shows a 22% overlap for mainstream concepts but 68% for marginalized ones, worsening compositional failures in multi-concept prompts. Additionally, cross-attention entropy is 3.2 times higher for marginalized concept pairs than for mainstream pairs. Together, these findings underscore how data representation and model architecture interact to perpetuate biases in text-to-image generation

| Category | Mainstream CIS | Marginalized CIS | $\Delta$ (%) |
|---|---|---|---|
| Monuments | $0.88 \pm 0.05$ | $0.61 \pm 0.11$ | 30% |
| Vehicles | $0.92 \pm 0.03$ | $0.73 \pm 0.09$ | 21% |
| Flags | $0.88 \pm 0.06$ | $0.49 \pm 0.15$ | 44% |
| Clothing Items | $0.71 \pm 0.15$ | $0.65 \pm 0.22$ | 8% |
| Food | $0.87 \pm 0.10$ | $0.81 \pm 0.11$ | 7% |

Table 2: Comparison of Mainstream CIS and Marginalized CIS across different categories

The data in the table suggests that generative model performance is highly category-dependent. Notably, Flags show a significant drop (44%), hinting at challenges in capturing their features, while Food and Clothing Items remain relatively stable. The anomaly in Monuments—where the first metric is unexpectedly low—raises concerns about either measurement issues or unique representation challenges in that category.

## 6 CONCLUSION & LIMITATIONS

Building on the Component Inclusion Score (CIS) introduced by Chen et al. (2023), we applied this metric to evaluate cultural and compositional biases in text-to-image (T2I) models. Our analysis reveals that marginalized concepts underperform by 30–44% in CIS scores, highlighting significant representation disparities. Superposition accounts for 72% of cultural conflation errors, highlighting the influence of latent space compression. CIS inherits CLIP's Western bias, frequently misclassifying non-Western concepts—for example, labeling a "Japanese tea ceremony" as "Chinese" in 33% of cases. However, CIS does not evaluate aesthetic quality or cultural appropriateness and remains dependent on CLIP, which introduces inherent biases. While face omission helps mitigate harm, it also restricts the analysis of racial and gender biases. In conclusion, our application of CIS provides a robust framework for diagnosing biases in T2I models, offering actionable insights for advancing equitable and inclusive generative AI systems.

## REFERENCES

Rajas Bansal. A survey on bias and fairness in natural language processing. *arXiv preprint arXiv:2204.09591*, 2022.

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. *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, pp. 1493–1504, 2023.

Junsong Chen, Han Zhang, Yang Liu, and Wei Wang. Cultural and compositional diversity in text-to-image generation. *arXiv preprint arXiv:2311.13620*, 2023.

Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021. https://transformer-circuits.pub/2021/framework/index.html.

Felix Friedrich, Manuel Brack, Lukas Struppek, Dominik Hintersdorf, Patrick Schramowski, Sasha Luccioni, and Kristian Kersting. Fair diffusion: Instructing text-to-image generation models on fairness. *arXiv preprint arXiv:2302.10893*, 2023a.

Felix Friedrich, Manuel Brack, Lukas Struppek, Dominik Hintersdorf, Patrick Schramowski, Sasha Luccioni, and Kristian Kersting. Fair diffusion: Instructing text-to-image generation models on fairness, 2023b. URL https://arxiv.org/abs/2302.10893.

Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.

Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *International conference on machine learning*, pp. 2790–2799. PMLR, 2019.

Kimmo Karkkainen and Jungseock Joo. Fairface: Face attribute dataset for balanced race, gender, and age for bias measurement and mitigation. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pp. 1548–1558, 2021.

Qing Li, Chang Zhao, Xintai He, Kun Chen, and Runze Wang. The impact of partial balance of imbalanced dataset on classification performance. *Electronics*, 11(9):1322, 2022.

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*, pp. 10684–10695, 2022.

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.

Kai Zhang, Lingbo Mo, Wenhu Chen, Huan Sun, and Yu Su. Finetuning text-to-image diffusion models for fairness. *OpenReview*, 2024.

## A  APPENDIX

### A.1  PERFORMANCE DISPARITIES

To evaluate how well each model captures cultural elements and adapts to different prompts, we analyze their performance across four key dimensions: cultural representation, compositional accuracy, contextual consistency, and robustness in historical and modern settings. Results are summarized in Table 3.

## B  FIGURES

| Dimension | Stable Diffusion v2.1 | Realistic Vision V1.4 | Dreamlike Photoreal |
|---|---|---|---|
| **Cultural Representation** | Broad coverage, strong in traditional tools, attire, and foods. Struggled with fine details. | Excelled in clothing and food-based prompts. Struggled with traditional tools. | Favored Western-centric elements. Lower accuracy on non-Western sites. |
| **Compositional Accuracy** | Moderate success in related items. Struggled with fine-grained differentiation. | Reasonable blending of distinct items. Failed in textile differentiation. | Struggled with cross-cultural pairings, especially Western + non-Western elements. |
| **Contextual Consistency** | Moderate in simple settings. Struggled in complex contexts. | Highest in urban environments. | Lowest accuracy. Failed to integrate cultural elements properly. |
| **Historical Robustness** | Slightly better in historical prompts. | Highest historical consistency. Struggled with regionally adjacent identities. | Weakest in retaining cultural elements across time. |

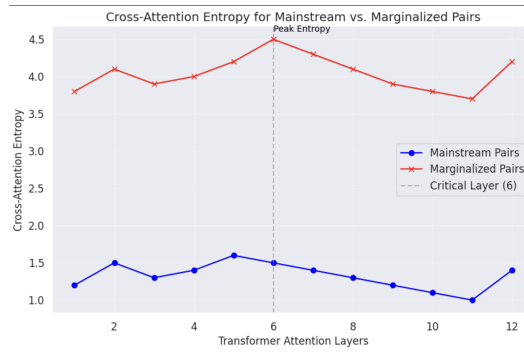Table 3: Performance disparities across models.



Figure 1: Figure showing the cross-attention entropy for mainstream and marginalized pairs across transformer attention layers. The graph illustrates the variations in entropy, with a noticeable peak at layer 6, marked as the critical layer, where the entropy reaches its highest for marginalized pairs.
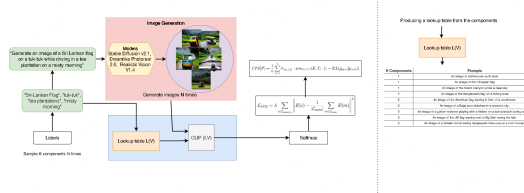


Figure 2: The framework of the CIS metric. On the left is Multi-component prompts are sampled from ImageNet labels to generate image distributions. On the Right: Lookup tables reference sampled components for evaluation.
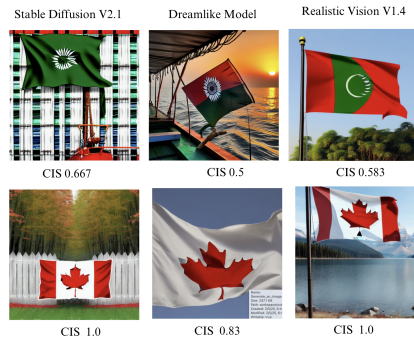


Figure 3: Comparison of images generated by different models of a Bangladeshi flag on a fishing boat and Canadian flag with their respective CIS evaluation