

Uncovering Cultural Representation Disparities in Vision-Language Models

Anonymous EMNLP submission

Abstract

Vision-Language Models (VLMs) have demonstrated impressive capabilities across a range of tasks, yet concerns about their potential biases exist. This work investigates the extent to which prominent VLMs exhibit cultural biases by evaluating their performance on an image-based country identification task at a country level. Utilising the geographically diverse Country211 (OpenAI, 2021) dataset, we probe several VLMs under various prompting strategies: open-ended questions, multiple-choice questions (MCQs), including challenging setups like multilingual and adversarial settings. Our analysis aims to uncover disparities in model accuracy across different countries and question formats, providing insights into how training data distribution and evaluation methodologies might influence cultural biases in VLMs. The findings highlight significant variations in performance, suggesting that while VLMs possess considerable visual understanding, they inherit biases from their pre-training data and scale that impact their ability to generalize uniformly across diverse global contexts.

1 Introduction

VLMs have rapidly advanced, demonstrating exceptional capabilities in integrating visual and textual information for a wide array of tasks, from image captioning to visual question answering (Liu et al., 2024; Alayrac et al., 2022; Wang et al., 2024). These models are increasingly being deployed in diverse applications, impacting areas such as education, healthcare, and public services globally (Zhang et al., 2024).

However, as their influence grows, so do concerns regarding their potential to perpetuate and

even amplify societal biases present in their training data (Zhao et al., 2017; Zhou et al., 2022; Weng et al., 2024). Cultural and geographical biases are of particular concern because they can lead to unequal performance and representation across different populations and regions of the world (AlKhamissi et al., 2024; Manvi et al., 2024). Defining "culture" is inherently complex, encompassing a broad spectrum of social norms, values, practices, languages, and historical contexts that shape the lived experiences of individuals and communities (Kroeber et al., 1985). Establishing culture in computational settings presents a persistent challenge due to its multifaceted and dynamic nature. Empirical studies employ tractable proxies such as demographic or geographic proxies to enable systematic analysis (Adilazuarda et al., 2024; Yadav et al., 2025). While nation-level aggregation can mask sub-national heterogeneity, prior work in human-computer interaction and cultural analytics has demonstrated that country labels often serve as a practical proxy for coarse-grained cultural signals when large-scale analyses are required (Obradovich et al., 2022).

In order to quantify cultural disparities in VLMs, we adopt image-based country identification as a concrete proxy task in which a model must infer an image's country of origin solely from visual cues, while also providing a justification. Prior work has shown that geolocation tasks reveal representational imbalances in visual models, as performance often correlates with the prevalence of training data from different regions (Pouget et al., 2024).

The main contributions of this paper are:

1. We introduce a scalable framework to evaluate cultural biases in VLMs using an image-based country identification task over 211 countries,

* Equal Contribution.

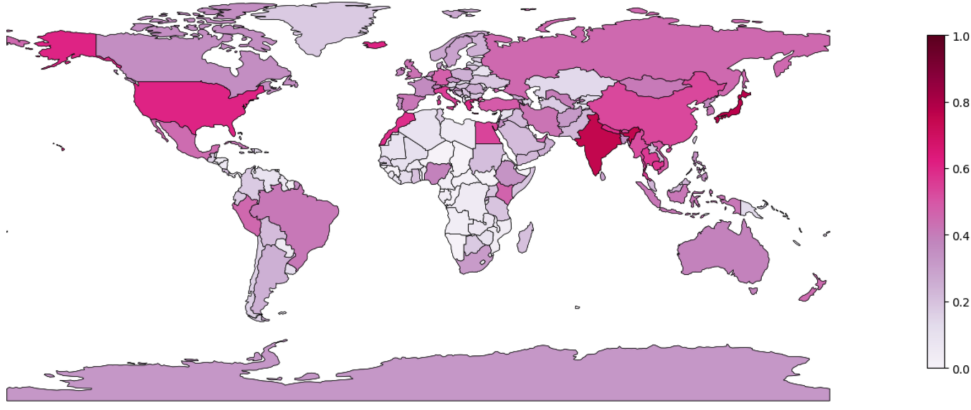


Figure 1: Visualization of the average country-wise recognition accuracy across the VLMs studied in this paper. VLMs perform well at recognizing images from North American and Western European countries, but there are clear disparities in performance for African and Central American countries.

leveraging the geographically diverse and balanced Country211 dataset.

2. We systematically probe VLMs under varied settings—open-ended and multiple-choice questions (MCQs) with both random and culturally similar distractors—alongside multilingual prompts in five languages, to capture nuanced cultural and linguistic disparities.¹
3. We examine model robustness to image perturbations and analyse performance across nine image categories (e.g. architecture, landscape, food etc), revealing the influence of image content on cultural bias.
4. Our findings show that VLM biases do not consistently favour Western countries; instead, biases often reflect over representation of certain popular countries (e.g., India, USA) in the training data², suggesting a more complex bias landscape.

2 Related Works

Recent work has increasingly explored the socio-cultural dimensions of Large Language Models (LLMs), including how they encode, express, and respond to culturally specific knowledge. Studies have examined value alignment (Choenni

and Shutova, 2024), moral reasoning across languages (Agarwal et al., 2024), and cultural persona (AlKhamissi et al., 2024), while also uncovering strong Western biases in model outputs (Naous et al., 2024) which risk marginalizing cultural diversity if deployed in real world. There have also been efforts to address these concerns, like prompting based on ethnographic fieldwork (AlKhamissi et al., 2024) and fine-tuning culture-specific LLMs (Li et al., 2024a). Similar studies have been extended for Vision Language Models (VLMs) starting from (Liu et al., 2021) over cultural aspects, but in a weaker capacity (Nwatu et al., 2023) showed that CLIP (Radford et al., 2021) struggled in data for poor socio-economic groups worldwide in the Dollar Street dataset (Gaviria Rojas et al., 2022). State-of-the-art off-shelf VLMs score much higher on images depicting Western scenes than equivalent East-Asian scenes for every vision task, such as identification, question-answering, and art emotion classification (Ananthram et al., 2025). Similarly, (Liu et al., 2025; Yadav et al., 2025) reveals that VLMs show stronger performance in Western concepts and weaker results in African and Asian contexts. These findings align with the fact that large pretraining corpora are dominated by high-resource languages and regions. Of the samples that can be geo-located in the OpenImages dataset (Kuznetsova et al., 2020), 32% were from only the United States, and 60% came from only six Western countries (Shankar et al., 2017). Such imbalances translate into a “Western bias” in model behavior (de Vries et al., 2019).

Datasets & Benchmarks : To probe these biases, a growing body of work has constructed specialized

¹Due to cultural similarities, misclassification among similar countries is more likely than misclassification with an unrelated country. MCQ with random and similar distractors tested the VLMs in both scenarios as to whether misclassification would occur when all distractors are neither neighboring nor similar countries

²For deliberately under specified inputs without country names, the generated images most reflect the surroundings of the United States followed by India. (Basu et al., 2023)

Prior Work	Eval Method	Multilingual?	Adversarial?	Categories	Total Sample Count	Domain
CulturalVQA (Nayak et al., 2024)	Open-Ended	No	No	11 Countries	2,328	5 Categories
WorldCuisines (Winata et al., 2025)	Both	Yes (30 languages)	Yes	189 Countries	6,045	Only Food
Food-500 CAP (Ma et al., 2023)	Open-Ended	No	Yes	7 Regions	24,700	Only Food
MOSAIC-1.5k (Burda-Lassen et al., 2025)	Open-Ended	No	No	N/A	1,500	3 Categories
See It From My Perspective (Ananthram et al., 2025)	Open-Ended	Yes (2 languages)	No	2 Regions	38,479	4 Categories
CVQA (Romero et al., 2024)	MCQ	Yes (31 languages)	Yes	39 Countries	5,239	10 Categories
GIMMICK (Schneider et al., 2025)	Both(MCQ)**	No	No	144 Countries	7,239(1,741)**	-
Ours	Both	Yes (5 languages)	Yes	211 Countries	63,300	9 Categories

Table 1: Overview of prior datasets used in cultural recognition experiments.

**The vales in brackets indicate the features in the Country recognition task subset, while the first values indicate that of the whole dataset.



Figure 2: Examples of the Country211 dataset, alongside automatically-predicted categories for each image, showcasing the visual diversity of the examples to be classified.

datasets and benchmarks with cross-cultural content, such as MOSAIC-1.5k (Burda-Lassen et al., 2025), CULTURAL-VQA (Nayak et al., 2024), and GlobalRG (Bhatia et al., 2024). Many works also opt for probing specific aspects of culture, such as food (Li et al., 2024b), race (tse Huang et al., 2025), art (Mohamed et al., 2024), etc., instead of providing an overall view for bias study. (Winata et al., 2025) introduced WorldCuisines for Food Vision Question Answering and country identification and found that VLMs often fail on adversarially misleading contexts or less-common cuisines. (Ma et al., 2023) introduced the Food-500 CAP dataset and observed that most models exhibited geographical culinary biases. Several studies have also treated country-of-origin or geolocation as a proxy for cultural provenance. WorldCuisines includes a country identification task to reveal failures on uncommon or misleading contexts (Winata et al., 2025), and Food-500 CAP finds systematic mismatches between predicted and actual countries of culinary images (Ma et al., 2023). Even in datasets like Dollar Street (Gaviria Rojas et al., 2022) or OpenImages (Kuznetsova et al., 2020), ge-

ographic metadata has been used to analyze representational imbalances across regions (Nwatu et al., 2023; Shankar et al., 2017), demonstrating that country-level annotations provide a practical signal for probing cultural and geographic bias in VLMs.

Impact of Evaluation: The format of evaluation also impacts bias measurement. Many of the above benchmarks use multiple-choice or binary questions, which can mask a model’s true understanding. Since language choice can influence bias, benchmarks are often performed across multiple languages. (Romero et al., 2024) showed that the performance of LLaVA-1.5-7B dropped by 19.6% when prompted without multiple choices for CVQA. Models also showed lower performance when prompted in native language of the image’s country of origin. However, (Ananthram et al., 2025) observed that prompting in a culturally closer language can reduce Western bias in some VLMs. It was also observed that people of different cultures are capable of differently capable of describing what they see in an image (van Miltenburg et al., 2017). We build on these insights by comparing open-ended vs. multiple-choice prompts (including

“hard” questions with challenging distractors) and by evaluating in both English and native languages, to see how the prompting strategy affects cultural bias in VLMs.

3 Dataset Used

The primary dataset used for the experiments is the Country211 (Radford et al., 2021) dataset which was a subset of images from YFCC100M (Thomee et al., 2016) having GPS coordinates associated with them. The images cover several domains including but not limited to - exterior architecture, interior architecture, landscape (vegetation, nature, skyview), people’s appearance, attires, scripts, texts, posters, etc. The GPS coordinates associated with the images were then used to map them to individual countries. ISO-3166³ codes representing each country were used as labels for each image. ISO labels were used for consistency as country names used by the VLMs were not deterministic i.e Britain was also used simultaneously in place of Great Britain or UK or its constituents, proving the list of tags and corresponding country names led to the models responding consistently with no observable difference in performance. For our experiments, we utilized this dataset, which consists of 21.1 K images, i.e 100 images each from 211 countries.

Key Differences: Existing benchmarks highlight cultural blind spots in VLMs, but they generally either cover fewer categories or countries or are restricted to specialized domains. Our work differs by using an image-based country-identification task over 211 countries, providing much broader geographic coverage and adversarial probing. Further, the datasets utilized in the prior works utilize a images that might be easier to classify, including but not limited to close up shots of food items, popular monuments being the primary object in an image etc. The dataset we utilized introduces a lot of noise and randomness in a majority of images as seen in Figure 2⁴. For instance, the examples shown from UK, India and Egypt might be

³https://en.wikipedia.org/wiki/List_of_ISO_3166_country_codes

⁴The images were part of OpenAI’s YFCC100M and come with pre-verified country labels. Although some samples might be difficult to classify even for a native, The primary goal was to uncover cultural biases using the features the VLMs could probably misclassify it with a culturally similar or neighboring country, but frequently misclassify it with a very dissimilar country.

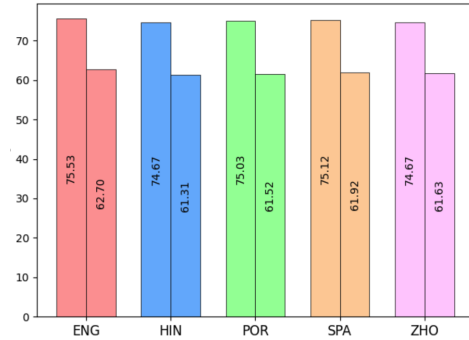


Figure 3: Model-wise averaged accuracy when varying the prompt language or selection of MCQA alternatives (left: random; right: similar). Performance is consistent across conditions.

easy to classify, but the examples from Afghanistan and Kuwait require grasping certain features and their associated knowledge i.e the headgear pattern of the Kuwait image and how it is different from other countries in the region. Th example from Afghanistan requires noticing the afghan flag, while the appearance of the person in the left may try to mislead the VLM due to an appearance of different ethnicity.

4 Experiments

Prompt Variations: We probed each VLM under three complementary prompting paradigms.

1. open-ended questions
2. MCQs (with random distractors)
3. MCQs (with similar distractors)

Image perturbations: Open-ended experiments were done with these adversarial changes:

1. Rotation by 90° clockwise,
2. Rotation by 90° anti-clockwise
3. Flipping the image
4. Gray-scaling the image

However, open ended experiments pose challenges for objective scoring due to semantic variability. Second, multiple-choice questions (MCQs) with random distractors yield correctness metrics yet may understate subtle biases if distractors are easily ruled out. Third, challenging MCQs with similar distractors force models to discriminate between culturally proximate options, thus exposing fine-grained bias patterns. The MCQs are designed as part of discriminative probing and to assess the disparity in the model’s

cultural knowledge.

Linguistic Variations : We further extend discriminative proving to a multilingual setting, prompting models in five languages : (English, Hindi, Chinese, Portuguese, Spanish) to assess the intersection of cultural and linguistic biases.

Model Variations : A diverse set of VLMs were tested including both proprietary and open-weight models of varying sizes: Gemini-2.5-Flash, Gemma-3-27B (Team et al., 2025), Aya-Vision-8B, Aya-Vision-32B (Dash et al., 2025), GPT-4o-Mini (OpenAI et al., 2024), (etal, 2025).

The experiments being repeated with each permutation of features lead to a total of 168.8 K samples tested. Inference was done in JSON format with the default hyperparameters for each of the models tested through Cohere⁵ and OpenRouter’s API⁶. More on the JSON formatting and prompts used can be found in Appendix D.

4.1 Open-Ended Evaluation

For the open-ended experiments, we asked each model to provide information on 4 areas: (1) name of the country, (2) country selection rationale in a few sentences, (3) a score from 0 to 100 representing the confidence in the classification, (4) and up to 6 features from the image as a list that had influence in the decision. The accuracies of each country obtained using each of the VLMs used can be seen in Figure 17. The accuracies of many countries were far lower especially in Eastern Europe, South America, Africa and Central Asia. This gap between country level accuracies was far higher in open ended experiments compared to the multiple-choice experiments .

4.2 Evaluation through random distractors in multiple-choice questions

For these experiments, we asked each model to provide information on 4 areas: (1) name of the country, (2) label of the chosen country from the choices provided (3) country selection rationale in a few sentences, and (4) a score from 0 to 100 representing the confidence in the classification. For these experiments, 4 countries were chosen at random from among the other 210 countries for each sample as distractors. The order of options were then shuffled such that the distribution of correct

answer’s location is made uniform. Compared to other settings, this setting led to the highest average of accuracies obtained due to the clearly contrasting nature of distractors used. However, many central African nations still face a recognition bias likely due to low representation in training data. This was observed across all VLMs that were tested as seen in Figure 18.

4.3 Evaluation through similar distractors in multiple-choice questions

Similar to the prior experiments with MCQs using random distractors, in this setting use similar nations as distractors. These were chosen from among the bordering nations. Any countries with high similarity in culture if any were added manually. (Ex : Spain -> Mexico). This led to the average of accuracies dropping considerably due to the challenging nature of the options presented to the models. However, the drops were observed for only a few countries where choosing similar distractors led to these countries’ images being classified as belonging to one of their popular neighbors. This can be observed in Figure 18 and Figure 19.

5 Results

The results for experimental setting over countries of each region can be seen in Table 3. Additionally Table 2 demonstrates the statistical significance of the results which presents Pearson’s χ^2 and Cochran’s Q test results for every evaluation condition and its subcategories. The χ^2 p-values of 0.0 (effectively underflowing) reveal that the six VLMs differ in accuracy in a highly significant way, confirming that VLM biases drive overall performance differences. In contrast, the Cochran’s Q p-values of 1.0 indicate no significant variation among sub-conditions (whether comparing original, rotated, or grayscale images, or across the five languages), showing that these perturbations do not meaningfully alter each VLM’s overall accuracy, but do cause changes in country wise accuracy as seen in Figure 14 and Figure 15.

5.1 Effect of Language of Inputs on Results

The average of country level accuracies compared to each language as input can be seen in Figure 3. The language used for inputs had a very little effect i.e <2% for all languages. But at a country level, most countries remained unaffected by language of the prompt to a large extent with change

⁵<https://docs.cohere.com/cohere-documentation>

⁶<https://openrouter.ai/docs/quickstart>

Condition	χ^2	p	Q	p _Q
Open-ended				
Original	5851.81	~ 0.0	-0.727	~ 1.0
Rotated	4666.49	~ 0.0	-0.952	~ 1.0
Greyscale	4280.08	~ 0.0	-0.725	~ 1.0
MCQ-S				
Overall	26274.855	~ 0.0	-0.151	~ 1.0
ENG	5232.00	~ 0.0	-0.145	~ 1.0
HIN	5197.50	~ 0.0	-0.153	~ 1.0
POR	5309.00	~ 0.0	-0.155	~ 1.0
SPA	5256.72	~ 0.0	-0.151	~ 1.0
ZHO	5345.21	~ 0.0	-0.155	~ 1.0
MCQ-R				
Overall	23855.61	~ 0.0	-0.074	~ 1.0
ENG	4835.90	~ 0.0	-0.073	~ 1.0
HIN	4584.62	~ 0.0	-0.073	~ 1.0
POR	4873.45	~ 0.0	-0.076	~ 1.0
SPA	4782.98	~ 0.0	-0.074	~ 1.0
ZHO	4801.08	~ 0.0	-0.076	~ 1.0

Table 2: Statistical tests for each evaluation condition and subcategory : Open ended, MCQ-Random (MCQ-R), MCQ-Similar (MCQ-S)

in accuracy less than 0.1%. The only cases with a noticeable change in accuracy are some but not all of the countries that speak the target language predominantly. For example, Changing the input language from English to Spanish improved accuracy for Spain but the change over Latin-American countries was negligible. Similarly, while switching to Portuguese had improved the accuracy for Brazil, it led to a drop in accuracy for Portugal. Overall, the input language improves performance for some countries primarily associated with the language used. The results also partially contradict prior findings that prompting in culturally similar languages reduces western bias (Ananthram et al., 2025).

5.2 Effect of Image Perturbations on Results

Figure 4 and Figure 5 display the changes in accuracy observed due to gray-scaling and rotating the images compared to the original images. Input image perturbations can have a large impact on the country-level biases in VLMs. Further, It can be assumed that the VLMs tested are not robust enough towards image perturbations, with each country being effected at a different scale between each model/perturbation. The overall averages can also be seen in Figure 8, Figure 9 and Figure 10

Region	MCQA		
	Open-Ended	Similar	Random
North America	41.9	73.7	80.2
Central America	11.1	69.7	68.0
Caribbean	13.6	50.5	71.4
South America	20.4	70.9	68.7
Oceania	19.0	57.5	68.9
Western Europe	30.9	57.9	77.5
Northern Europe	25.3	60.6	79.4
Eastern Europe	26.6	53.4	75.9
Middle East	29.3	68.4	77.1
Central Asia	26.7	53.5	78.1
East Asia	43.6	71.6	83.8
Southeast Asia	41.7	67.5	81.7
South Asia	49.1	69.0	85.5
North Africa	31.9	54.3	78.9
Central Africa	11.8	57.0	68.2
Southern Africa	20.4	74.2	74.2
Overall	27.7	63.1	76.1

Table 3: Region-wise averaged accuracy across models. There are consistent disparities in performance across different regions, regardless of the prompting method.

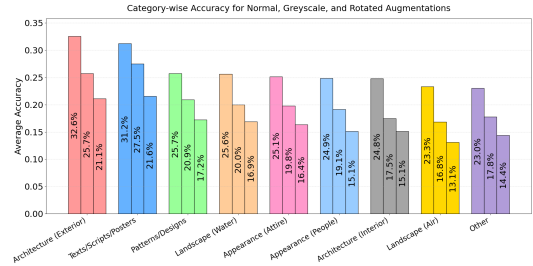


Figure 4: Model-wise averaged accuracy across the nine image categories, as a function of the image perturbations. There is a clear trend of models performing better with the original images (left), compared to the grayscale images (middle), or rotated images (right).

respectively.

Figure 18 shows how perturbations affect model performance across different semantic image categories. For all nine categories, models perform best on original (unaltered) images, with decreasing accuracy for gray-scaled and worse for rotated versions. Categories like exterior architecture, text/scripts/posters, and attire/patterns are especially impacted by perturbations. We hypothesize that it is likely because they contain fine-grained, orientation-sensitive, or highly color-dependent details.

We also look at geographical disparities of these changes in orientation in Figure 14 and Figure 15.

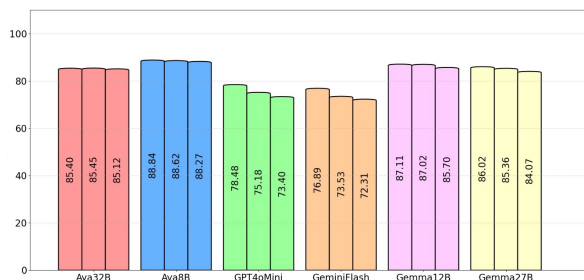


Figure 5: Average model confidence, given the original images (left), grayscale images (middle), and rotated images (right). GPT4o, GeminiFlash, and Gemma27B are most sensitive to image perturbations.

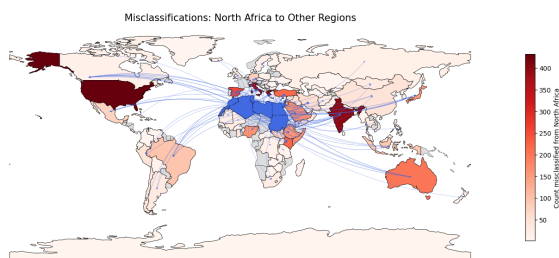


Figure 6: Mis-classification map for North African countries. There is a clear trend of models predicting USA, India, Australia, or geographically close countries in Europe and the Middle East.

We observe the disparity in model robustness also emerges clearly. For example, models such as Aya Vision 32B, GPT4o-mini and Gemini 3 12B show very different sensitivity across both a) perturbations and b) regions which were affected. We hypothesise that architectural and training differences might be influencing how models process image orientation and color. While gray-scaling may reduce performance due to the loss of visual detail or color-dependent cues, rotation disrupts spatial reasoning and object orientation, which are critical for geographic or cultural recognition.

These findings highlight the importance of evaluating model performance under realistic image distortions, especially for applications where images may not be clean or consistently formatted as image characteristics can vary widely.

5.3 Effect of Input Variations on Confidence

Despite the drop in Overall accuracy by all of the tested models due to either of the image perturbations, the confidence of the open-weight models didn't have a significant change while the proprietary models displayed a visible drop in confidence compared to the original images. Compared to rotation of images, Gray-scaling had a larger impact

on the response accuracies. The average confidence of each VLM with each adversarial setting compared to the original can be seen in Figure 5. The closed-weights models exhibited a drop in confidence when a rotated or grayscale image was provided than the corresponding originals, but this wasn't the case with open-weight models we tested.

5.4 Image Feature categories VS accuracy

Apart from the experiments, the original 21.1 K images were also labeled multi-way based on the key features they contain using larger VLMs like Gemini-2.5-Pro, o4-mini, Grok-2-Vision. Later a majority vote of each label was considered. The quality was later manually verified over a subset by multiple people⁷. We have used 9 sub-categories for this categorization. The descriptions of each of these categories can be seen in Table 4. A large variance was observed between each feature category and the country level accuracies obtained. Additionally there was also a large variation between how accuracy was affected for each country/feature based on model/perturbation used. This can be also be seen in Figure 13. The extent to which each category's images were recognized by VLMs can be seen in Figure 4. External architecture and native language texts' presence in the background helped the VLMs recognize the culture better compared to the other features.

5.5 Distribution of Predicted countries

The distribution of responses in an open ended approach can be seen in Figure 7. The output distributions varied largely among models, even those within the same family (i.e between Gemma-3-27B, Gemma-3-12B and Aya-vision-32B, Aya-vision-8B). The results obtained contradict the usual assumption about western biases in generative models, and was observed over a few nations with likely high training data proportion.

Notably, all models consistently overpredict certain countries—particularly the USA, India, and Brazil—regardless of actual ground truth. We hypothesize that these countries are likely overrepresented in the models' pretraining data or benefit from more visually distinctive cues. Biases seem to cluster around a few highly represented or visually salient countries rather than reflecting broader

⁷Feature category labels were verified on a subset of 10% samples equally distributed over all countries, with 2 people verifying labels, in cases with no consensus between the two, the third annotator was used.

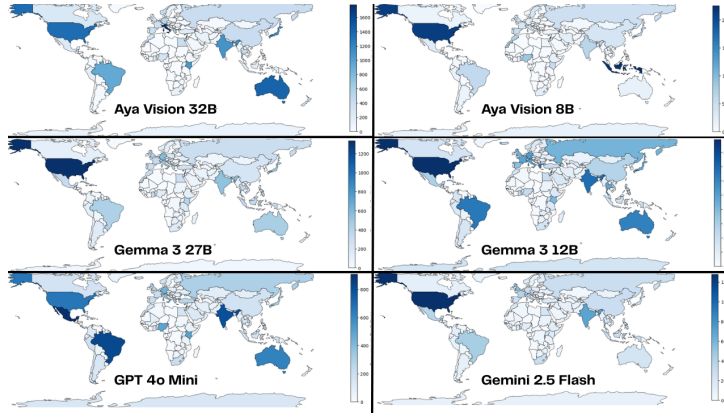


Figure 7: Country-wise response distribution in the open-ended prompt format. There is a consistent trend of models predicting USA, but otherwise, no clear bias towards predicting Western countries.

geopolitical landscape.

These results show that model predictions are likely highly influenced by data availability and image characteristics rather than a generic global bias. It also underscores the need for better interpretability regarding the geographic composition of VLM training datasets to fully understand such biases.

5.6 Misclassification Analysis

The mapping of misclassification of samples was not limited to similar or neighboring nations. This can be observed in Figure 20 to Figure 34. These misclassifications varied by each individual feature and provide a better fine-grained insights of cultural biases. For instance, Apart from neighboring / similar countries, most images from Africa and rural regions of South America were classified as India. A specific example is shown in Figure 6 where out of the 600 images (100 * 6 models), roughly 80-120 belong to this category for most countries, while many countries had most of their misclassified as originating from India.

6 Discussion

Our study presents a comprehensive analysis of cultural biases in Vision-Language Models (VLMs) using a geographically balanced dataset across 211 countries. We evaluated popular models across multiple prompting strategies, e.g. open-ended, multiple-choice (random and similar distractors), and multilingual settings. Open-ended formats showed the greatest disparity in country-level accuracy, particularly in underrepresented regions such as Central Africa and parts of South America. The use of culturally similar distractors proved to be

the most effective in revealing fine-grained errors, highlighting limitations in models’ cultural discrimination abilities.

We further assessed the models’ robustness to image perturbations like gray-scaling and rotation. While gray-scaling affected only a few specific countries, rotation led to a broad and uniform drop in performance, confirming that VLMs rely heavily on image orientation. We further observed that performance also varied by semantic image content—categories like architecture, textual cues, and attire were more predictive of cultural origin, especially in unaltered images. Language variation in prompts had minimal impact on average accuracy, though countries closely tied to the input language (e.g., Spain with Spanish, Brazil with Portuguese) showed slight gains. However, this trend was inconsistent and did not generalize across all culturally linked regions.

Finally, our misclassification analysis shows that models frequently confuse images from low-resource or visually ambiguous countries with a few dominant nations, reinforcing the role of training data bias.

7 Conclusion

Our findings show that biases are not uniformly Western but instead reflect over representation of certain countries in training data. Model performance varied across prompt types, languages, image features, and perturbations—highlighting limitations in robustness and cultural generalization. These results call for greater transparency in dataset composition and the need for more culturally inclusive evaluation methods to ensure fairer and more globally representative VLMs.

Limitations

Our study has a few important limitations to keep in mind. First, the use of country-level labels as a proxy for culture, while common for large-scale analysis, inherently overlooks intra-country cultural diversity and multicultural populations, potentially obscuring sub-national or regional nuances. The country labels used don't account for political complexities like disputed territories.

References

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

Utkarsh Agarwal, Kumar Tanmay, Aditi Khandelwal, and Monojit Choudhury. 2024. [Ethical reasoning and moral value alignment of LLMs depend on the language we prompt them in](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 6330–6340, Torino, Italia. ELRA and ICCL.

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob L Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikołaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. 2022. [Flamingo: a visual language model for few-shot learning](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 23716–23736. Curran Associates, Inc.

Badr AlKhamissi, Muhammad ElNokrashy, Mai Alkhamissi, and Mona Diab. 2024. [Investigating cultural alignment of large language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12404–12422, Bangkok, Thailand. Association for Computational Linguistics.

Amith Ananthram, Elias Stengel-Eskin, Mohit Bansal, and Kathleen McKeown. 2025. [See it from my perspective: How language affects cultural bias in image understanding](#). In *The Thirteenth International Conference on Learning Representations*.

Abhipsa Basu, R. Venkatesh Babu, and Danish Pruthi. 2023. [Inspecting the geographical representativeness of images from text-to-image models](#).

Mehar Bhatia, Sahithya Ravi, Aditya Chinchure, EunJeong Hwang, and Vered Shwartz. 2024. [From local concepts to universals: Evaluating the multi-cultural understanding of vision-language models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 6763–6782, Miami, Florida, USA. Association for Computational Linguistics.

Olena Burda-Lassen, Aman Chadha, Shashank Goswami, and Vinija Jain. 2025. [How culturally aware are vision-language models?](#) In *2025 IEEE 6th International Conference on Image Processing, Applications and Systems (IPAS)*, volume CFP2540Z-ART, pages 1–6.

Rochelle Choenni and Ekaterina Shutova. 2024. [Self-alignment: Improving alignment of cultural values in llms via in-context learning](#).

Saurabh Dash, Yiyang Nan, John Dang, Arash Ahmadian, Shivalika Singh, Madeline Smith, Bharat Venkitesh, Vlad Shmyhlo, Viraat Aryabumi, Walter Beller-Morales, Jeremy Pekmez, Jason Ozuzu, Pierre Richemond, Acyr Locatelli, Nick Frosst, Phil Blunsom, Aidan Gomez, Ivan Zhang, Marzieh Fadaee, Manoj Govindassamy, Sudip Roy, Matthias Gallé, Beyza Ermis, Ahmet Üstün, and Sara Hooker. 2025. [Aya vision: Advancing the frontier of multilingual multimodality](#).

Terrance de Vries, Ishan Misra, Changhan Wang, and Laurens van der Maaten. 2019. Does object recognition work for everyone? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*.

Gemini Team ... etal. 2025. [Gemini: A family of highly capable multimodal models](#).

William Gaviria Rojas, Sudnya Diamos, Keertan Kini, David Kanter, Vijay Janapa Reddi, and Cody Coleman. 2022. [The dollar street dataset: Images representing the geographic and socioeconomic diversity of the world](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 12979–12990. Curran Associates, Inc.

A. L. Kroeber, Wayne Untereiner, and Clyde Kluckhohn. 1985. *Culture: A critical review of concepts and definitions*. Vintage Books.

Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, Tom Duerig, and Vittorio Ferrari. 2020. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *IJCV*.

Cheng Li, Mengzhuo Chen, Jindong Wang, Sunayana Sitaram, and Xing Xie. 2024a. [Culturellm: Incorporating cultural differences into large language models](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 84799–84838. Curran Associates, Inc.

- Wenyan Li, Crystina Zhang, Jiaang Li, Qiwei Peng, Raphael Tang, Li Zhou, Weijia Zhang, Guimin Hu, Yifei Yuan, Anders Søgaard, Daniel Hershcovich, and Desmond Elliott. 2024b. [FoodieQA: A multi-modal dataset for fine-grained understanding of Chinese food culture](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 19077–19095, Miami, Florida, USA. Association for Computational Linguistics.
- Fangyu Liu, Emanuele Bugliarello, Edoardo Maria Ponti, Siva Reddy, Nigel Collier, and Desmond Elliott. 2021. [Visually grounded reasoning across languages and cultures](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10467–10485, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 26296–26306.
- Shudong Liu, Yiqiao Jin, Cheng Li, Derek F. Wong, Qingsong Wen, Lichao Sun, Haipeng Chen, Xing Xie, and Jindong Wang. 2025. [Culturevlm: Characterizing and improving cultural understanding of vision-language models for over 100 countries](#).
- Zheng Ma, Mianzhi Pan, Wenhan Wu, Kanzhi Cheng, Jianbing Zhang, Shujian Huang, and Jiajun Chen. 2023. Food-500 cap: A fine-grained food caption benchmark for evaluating vision-language models. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 5674–5685.
- Rohin Manvi, Samar Khanna, Marshall Burke, David B. Lobell, and Stefano Ermon. 2024. [Large language models are geographically biased](#). In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 34654–34669. PMLR.
- Youssef Mohamed, Runjia Li, Ibrahim Said Ahmad, Kilichbek Haydarov, Philip Torr, Kenneth Church, and Mohamed Elhoseiny. 2024. [No culture left behind: ArtELingo-28, a benchmark of WikiArt with captions in 28 languages](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 20939–20962, Miami, Florida, USA. Association for Computational Linguistics.
- Tarek Naous, Michael J Ryan, Alan Ritter, and Wei Xu. 2024. [Having beer after prayer? measuring cultural bias in large language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16366–16393, Bangkok, Thailand. Association for Computational Linguistics.
- Shravan Nayak, Kanishk Jain, Rabiul Awal, Siva Reddy, Sjoerd Van Steenkiste, Lisa Anne Hendricks, Karolina Stanczak, and Aishwarya Agrawal. 2024. [Benchmarking vision language models for cultural understanding](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 5769–5790, Miami, Florida, USA. Association for Computational Linguistics.
- Joan Nwatu, Oana Ignat, and Rada Mihalcea. 2023. [Bridging the digital divide: Performance variation across socio-economic factors in vision-language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10686–10702, Singapore. Association for Computational Linguistics.
- Nick Obradovich, Ömer Özak, Ignacio Martín, Ignacio Ortuño-Ortín, Edmond Awad, Manuel Cebrián, Rubén Cuevas, Klaus Desmet, Iyad Rahwan, and Ángel Cuevas. 2022. Expanding the measurement of culture with a sample of two billion humans. *Journal of the Royal Society Interface*, 19(190):20220085.
- OpenAI. 2021. [Country211](#).
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan

762	Leike, Jade Leung, Daniel Levy, Chak Ming Li,	<i>Proceedings of Machine Learning Research</i> , pages	823
763	Rachel Lim, Molly Lin, Stephanie Lin, Mateusz	8748–8763. PMLR.	824
764	Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue,		
765	Anna Makanju, Kim Malfacini, Sam Manning, Todor	David Romero, Chenyang Lyu, Haryo Akbarianto Wi-	825
766	Markov, Yaniv Markovski, Bianca Martin, Katie	bowo, Teresa Lynn, Injy Hamed, Aditya Nanda	826
767	Mayer, Andrew Mayne, Bob McGrew, Scott Mayer	Kishore, Aishik Mandal, Alina Dragonetti, Artem	827
768	McKinney, Christine McLeavey, Paul McMillan,	Abzaliev, Atnafu Lambebo Tonja, Bontu Fufa	828
769	Jake McNeil, David Medina, Aalok Mehta, Jacob	Balcha, Chenxi Whitehouse, Christian Salamea,	829
770	Menick, Luke Metz, Andrey Mishchenko, Pamela	Dan John Velasco, David Ifeoluwa Adelani, David	830
771	Mishkin, Vinnie Monaco, Evan Morikawa, Daniel	Le Meur, Emilio Villa-Cueva, Fajri Koto, Fauzan	831
772	Mossing, Tong Mu, Mira Murati, Oleg Murk, David	Farooqui, Frederico Belcavello, Ganzorig Batnasan,	832
773	Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak,	Gisela Vallejo, Grainne Caulfield, Guido Ivetta,	833
774	Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh,	Haiyue Song, Henok Biadgign Ademteu, Hernán	834
775	Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex	Maina, Holy Lovenia, Israel Abebe Azime, Jan	835
776	Paino, Joe Palermo, Ashley Pantuliano, Giambat-	Christian Blaise Cruz, Jay Gala, Jiahui Geng,	836
777	tista Parascandolo, Joel Parish, Emy Parparita, Alex	Jesus-German Ortiz-Barajas, Jinheon Baek, Joce-	837
778	Passos, Mikhail Pavlov, Andrew Peng, Adam Perel-	lyn Dunstan, Laura Alonso Alemany, Kumaran-	838
779	man, Filipe de Avila Belbute Peres, Michael Petrov,	age Ravindu Yasas Nagasinghe, Luciana Benotti,	839
780	Henrique Ponde de Oliveira Pinto, Michael, Poko-	Luis Fernando D’ Haro, Marcelo Viridiano, Mar-	840
781	rny, Michelle Pokrass, Vitchyr H. Pong, Tolly Pow-	cos Estecha-Garitagoitia, Maria Camila Buitrago	841
782	ell, Alethea Power, Boris Power, Elizabeth Proehl,	Cabrera, Mario Rodríguez-Cantelar, Mélanie Jouis-	842
783	Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh,	teau, Mihail Mihaylov, Naome Etori, Mohamed Fa-	843
784	Cameron Raymond, Francis Real, Kendra Rimbach,	zli Mohamed Imam, Muhammad Farid Adilazuarda,	844
785	Carl Ross, Bob Rotsted, Henri Roussez, Nick Ry-	Munkhjargal Gochoo, Munkh-Erdene Otgonbold,	845
786	der, Mario Saltarelli, Ted Sanders, Shibani Santurkar,	Olivier Niyomugisha, Paula Mónica Silva, Pranjal	846
787	Girish Sastry, Heather Schmidt, David Schnurr, John	Chitale, Raj Dabre, Rendi Chevi, Ruochen Zhang,	847
788	Schulman, Daniel Selsam, Kyla Sheppard, Toki	Ryandito Diandaru, Samuel Cahyawijaya, Santiago	848
789	Sherbakov, Jessica Shieh, Sarah Shoker, Pranav	Góngora, Soyeong Jeong, Sukannya Purkayastha,	849
790	Shyam, Szymon Sidor, Eric Sigler, Maddie Simens,	Tatsuki Kuribayashi, Teresa Clifford, Thanmay	850
791	Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin	Jayakumar, Tiago Timponi Torrent, Toqeer Ehsan,	851
792	Sokolowsky, Yang Song, Natalie Staudacher, Fel-	Vladimir Araujo, Yova Kementchedjhieva, Zara	852
793	ipe Petroski Such, Natalie Summers, Ilya Sutskever,	Burzo, Zheng Wei Lim, Zheng Xin Yong, Oana Ignat,	853
794	Jie Tang, Nikolas Tezak, Madeleine B. Thompson,	Joan Nwatu, Rada Mihalcea, Tamar Solorio, and	854
795	Phil Tillet, Amin Tootoonchian, Elizabeth Tseng,	Alham Fikri Aji. 2024. Cvqa: Culturally-diverse mul-	855
796	Preston Tuggle, Nick Turley, Jerry Tworek, Juan Fe-	tilingual visual question answering benchmark . In	856
797	lipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,	<i>Advances in Neural Information Processing Systems</i> ,	857
798	Chelsea Voss, Carroll Wainwright, Justin Jay Wang,	volume 37, pages 11479–11505. Curran Associates,	858
799	Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei,	Inc.	859
800	CJ Weinmann, Akila Welihinda, Peter Welinder, Ji-		
801	ayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner,	Florian Schneider, Carolin Holtermann, Chris Biemann,	860
802	Clemens Winter, Samuel Wolrich, Hannah Wong,	and Anne Lauscher. 2025. GIMMICK: Globally	861
803	Lauren Workman, Sherwin Wu, Jeff Wu, Michael	inclusive multimodal multitask cultural knowledge	862
804	Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qim-	benchmarking . In <i>Findings of the Association for</i>	863
805	ing Yuan, Wojciech Zaremba, Rowan Zellers, Chong	<i>Computational Linguistics: ACL 2025</i> , pages 9605–	864
806	Zhang, Marvin Zhang, Shengjia Zhao, Tianhao	9668, Vienna, Austria. Association for Computa-	865
807	Zheng, Juntang Zhuang, William Zhuk, and Barret	tional Linguistics.	866
808	Zoph. 2024. Gpt-4 technical report .		
809	Angéline Pouget, Lucas Beyer, Emanuele Bugliarello,	Shreya Shankar, Yoni Halpern, Eric Breck, James At-	867
810	Xiao Wang, Andreas Peter Steiner, Xiaohua Zhai, and	wood, Jimbo Wilson, and D. Sculley. 2017. No clas-	868
811	Ibrahim Alabdulmohsin. 2024. No filter: Cultural	sification without representation: Assessing geodiver-	869
812	and socioeconomic diversity in contrastive vision-	sity issues in open data sets for the developing world.	870
813	language models . In <i>Advances in Neural Informa-</i>	In <i>NIPS 2017 workshop: Machine Learning for the</i>	871
814	<i>tion Processing Systems</i> , volume 37, pages 106474–	<i>Developing World</i> .	872
815	106496. Curran Associates, Inc.		
816	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya	Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya	873
817	Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-	Pathak, Nino Vieillard, Ramona Merhej, Sarah Per-	874
818	try, Amanda Askell, Pamela Mishkin, Jack Clark,	rin, Tatiana Matejovicova, Alexandre Ramé, Mor-	875
819	Gretchen Krueger, and Ilya Sutskever. 2021. Learn-	gane Rivière, Louis Rouillard, Thomas Mesnard, Ge-	876
820	ing transferable visual models from natural language	offrey Cideron, Jean bastien Grill, Sabela Ramos,	877
821	supervision . In <i>Proceedings of the 38th International</i>	Edouard Yvinec, Michelle Casbon, Etienne Pot, Ivo	878
822	<i>Conference on Machine Learning</i> , volume 139 of	Penchev, Gaël Liu, Francesco Visin, Kathleen Ke-	879
		nealy, Lucas Beyer, Xiaohai Zhai, Anton Tsitsulin,	880
		Robert Busa-Fekete, Alex Feng, Noveen Sachdeva,	881
		Benjamin Coleman, Yi Gao, Basil Mustafa, Iain	882

883	Barr, Emilio Parisotto, David Tian, Matan Eyal,	data in multimedia research. <i>Communications of the</i>	945
884	Colin Cherry, Jan-Thorsten Peter, Danila Sinopal-	<i>ACM</i> , 59(2):64–73.	946
885	nikov, Surya Bhupatiraju, Rishabh Agarwal, Mehran		
886	Kazemi, Dan Malkin, Ravin Kumar, David Vilar,	Jen tse Huang, Jiantong Qin, Jianping Zhang, Youliang	947
887	Idan Brusilovsky, Jiaming Luo, Andreas Steiner,	Yuan, Wenxuan Wang, and Jieyu Zhao. 2025. <i>Vis-</i>	948
888	Abe Friesen, Abhanshu Sharma, Abheesht Sharma,	<i>bias: Measuring explicit and implicit social biases in</i>	949
889	Adi Mayrav Gilady, Adrian Goedeckemeyer, Alaa	<i>vision language models</i> .	950
890	Saade, Alex Feng, Alexander Kolesnikov, Alexei		
891	Bendebury, Alvin Abdagic, Amit Vadi, András	Emiel van Miltenburg, Desmond Elliott, and Piek	951
892	György, André Susano Pinto, Anil Das, Ankur	Vossen. 2017. <i>Cross-linguistic differences and simi-</i>	952
893	Bapna, Antoine Miech, Antoine Yang, Antonia Pater-	<i>larities in image descriptions</i> . In <i>Proceedings of the</i>	953
894	son, Ashish Shenoy, Ayan Chakrabarti, Bilal Piot,	<i>10th International Conference on Natural Language</i>	954
895	Bo Wu, Bobak Shahriari, Bryce Petrini, Charlie	<i>Generation</i> , pages 21–30, Santiago de Compostela,	955
896	Chen, Charline Le Lan, Christopher A. Choquette-	Spain. Association for Computational Linguistics.	956
897	Choo, CJ Carey, Cormac Brick, Daniel Deutsch,		
898	Danielle Eisenbud, Dee Cattle, Derek Cheng, Dim-	Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi	957
899	itris Paparas, Divyashree Shivakumar Sreepathi-	Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang,	958
900	halli, Doug Reid, Dustin Tran, Dustin Zelle, Eric	Lei Zhao, Xixuan Song, Jiazheng Xu, Keqin Chen,	959
901	Noland, Erwin Huizenga, Eugene Kharitonov, Fred-	Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie	960
902	erick Liu, Gagik Amirhanyan, Glenn Cameron,	Tang. 2024. <i>Cogvlm: Visual expert for pretrained</i>	961
903	Hadi Hashemi, Hanna Klimczak-Plucińska, Har-	<i>language models</i> . In <i>Advances in Neural Informa-</i>	962
904	man Singh, Harsh Mehta, Harshal Tushar Lehri,	<i>tion Processing Systems</i> , volume 37, pages 121475–	963
905	Hussein Hazimeh, Ian Ballantyne, Idan Szpektor,	121499. Curran Associates, Inc.	964
906	Ivan Nardini, Jean Pouget-Abadie, Jetha Chan, Joe		
907	Stanton, John Wieting, Jonathan Lai, Jordi Orbay,	Zhaotian Weng, Zijun Gao, Jerone Andrews, and Jieyu	965
908	Joseph Fernandez, Josh Newlan, Ju yeong Ji, Jy-	Zhao. 2024. <i>Images speak louder than words: Un-</i>	966
909	otinder Singh, Kat Black, Kathy Yu, Kevin Hui, Ki-	<i>derstanding and mitigating bias in vision-language</i>	967
910	ran Vodrahalli, Klaus Greff, Linhai Qiu, Marcella	<i>model from a causal mediation perspective</i> . In <i>Pro-</i>	968
911	Valentine, Marina Coelho, Marvin Ritter, Matt Hoff-	<i>ceedings of the 2024 Conference on Empirical Meth-</i>	969
912	man, Matthew Watson, Mayank Chaturvedi, Michael	<i>ods in Natural Language Processing</i> , pages 15669–	970
913	Moynihan, Min Ma, Nabila Babar, Natasha Noy,	15680, Miami, Florida, USA. Association for Com-	971
914	Nathan Byrd, Nick Roy, Nikola Momchev, Nilay	putational Linguistics.	972
915	Chauhan, Noveen Sachdeva, Oskar Bunyan, Pankil		
916	Botarda, Paul Caron, Paul Kishan Rubenstein, Phil	Genta Indra Winata, Frederikus Hudi, Patrick Amadeus	973
917	Culliton, Philipp Schmid, Pier Giuseppe Sessa, Ping-	Irawan, David Anugraha, Rifki Afina Putri, Yutong	974
918	mei Xu, Piotr Stanczyk, Pouya Tafti, Rakesh Shiv-	Wang, Adam Nohejl, Ubaidillah Ariq Prathama, Ned-	975
919	anna, Renjie Wu, Renke Pan, Reza Rokni, Rob	jma Ousidhoum, Afifa Amriani, Anar Rzayev, Anir-	976
920	Willoughby, Rohith Vallu, Ryan Mullins, Sammy	ban Das, Ashmari Pramodya, Aulia Adila, Bryan	977
921	Jerome, Sara Smoot, Sertan Girgin, Shariq Iqbal,	Wilie, Candy Olivia Mawalim, Ching Lam Cheng,	978
922	Shashir Reddy, Shruti Sheth, Siim Pöder, Sijal Bhat-	Daud Abolade, Emmanuele Chersoni, Enrico San-	979
923	nagar, Sindhu Raghuram Panyam, Sivan Eiger, Susan	tus, Fariz Ikhwantri, Garry Kuwanto, Hanyang	980
924	Zhang, Tianqi Liu, Trevor Yacovone, Tyler Liechty,	Zhao, Haryo Akbarianto Wibowo, Holy Lovenia,	981
925	Uday Kalra, Utku Evci, Vedant Misra, Vincent Rose-	Jan Christian Blaise Cruz, Jan Wira Gotama Pu-	982
926	berry, Vlad Feinberg, Vlad Kolesnikov, Woohyun	tra, Junho Myung, Lucky Susanto, Maria Angel-	983
927	Han, Woosuk Kwon, Xi Chen, Yinlam Chow, Yuvein	ica Riera Machin, Marina Zhukova, Michael Anu-	984
928	Zhu, Zichuan Wei, Zoltan Egyed, Victor Cotruta,	graha, Muhammad Farid Adilazuarda, Natasha San-	985
929	Minh Giang, Phoebe Kirk, Anand Rao, Kat Black,	tosa, Peerat Limkonchotiawat, Raj Dabre, Rio Alexan-	986
930	Nabila Babar, Jessica Lo, Erica Moreira, Luiz Gus-	der Audino, Samuel Cahyawijaya, Shi-Xiong Zhang,	987
931	tavo Martins, Omar Sanseviero, Lucas Gonzalez,	Stephanie Yulia Salim, Yi Zhou, Yinxuan Gui,	988
932	Zach Gleicher, Tris Warkentin, Vahab Mirrokni, Evan	David Ifeoluwa Adelani, En-Shiun Annie Lee, Shogo	989
933	Senter, Eli Collins, Joelle Barral, Zoubin Ghahra-	Okada, Ayu Purwarianti, Alham Fikri Aji, Taro	990
934	mani, Raia Hadsell, Yossi Matias, D. Sculley, Slav	Watanabe, Derry Tanti Wijaya, Alice Oh, and Chong-	991
935	Petrov, Noah Fiedel, Noam Shazeer, Oriol Vinyals,	Wah Ngo. 2025. <i>Worldcuisines: A massive-scale</i>	992
936	Jeff Dean, Demis Hassabis, Koray Kavukcuoglu,	<i>benchmark for multilingual and multicultural visual</i>	993
937	Clement Farabet, Elena Buchatskaya, Jean-Baptiste	<i>question answering on global cuisines</i> .	994
938	Alayrac, Rohan Anil, Dmitry, Lepikhin, Sebastian		
939	Borgeaud, Olivier Bachem, Armand Joulin, Alek An-	Srishti Yadav, Zhi Zhang, Daniel Hershcovich, and Eka-	995
940	dreev, Cassidy Hardin, Robert Dadashi, and Léonard	terina Shutova. 2025. Beyond words: Exploring	996
941	Hussenot. 2025. <i>Gemma 3 technical report</i> .	cultural value sensitivity in multimodal models. In	997
		<i>Findings of the Association for Computational Lin-</i>	998
		<i>guistics: NAACL 2025</i> , pages 7592–7608.	999
942	Bart Thomee, David A Shamma, Gerald Friedland, Ben-		
943	jamin Elizalde, Karl Ni, Douglas Poland, Damian	Jingyi Zhang, Jiaying Huang, Sheng Jin, and Shijian Lu.	1000
944	Borth, and Li-Jia Li. 2016. Yfcc100m: The new	2024. <i>Vision-language models for vision tasks: A</i>	1001
		<i>survey</i> . <i>IEEE Transactions on Pattern Analysis and</i>	1002
		<i>Machine Intelligence</i> , 46(8):5625–5644.	1003

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2017. [Men also like shopping: Reducing gender bias amplification using corpus-level constraints](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2979–2989, Copenhagen, Denmark. Association for Computational Linguistics.

Kankan Zhou, Eason Lai, and Jing Jiang. 2022. [VL-StereoSet: A study of stereotypical bias in pre-trained vision-language models](#). In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 527–538, Online only. Association for Computational Linguistics.

A Overall Accuracies Before and After Image Perturbations

[Figure 8](#), [Figure 9](#), [Figure 10](#) display the accuracy obtained for each image perturbation used compared to the original through each of the VLMs tested.

B Overall Accuracy VS Models used : In each MCQ setting

[Figure 11](#), [Figure 12](#) display the accuracy obtained through each model in each MCQ experiment.

C Reproducibility

Inference was done through Cohere’s API for Aya-Vision-8B and Aya-Vision-32B through the default hyperparameters with a seed value of 1024. The rest of the models were used through OpenRouter’s API through the default hyper-parameters with a seed value of 1024. The experiments were repeated thrice and the overall accuracy varied between 1-1.2%, with some countries’ accuracy varying up to 1.5%. The costs associated with all experiments combined were 850\$ through OpenRouter and 250\$ Cohere API credits. The experiments were run on TPUs costing 0.35\$/hr with the costs reaching 60\$.

D Prompt Usage

The prompts used for each experiment setting can be seen below.

```
## OPEN ENDED
Respond in the following format :
{
  "Country": {
    "type": "string",
```

```
    "description": "Just Country Name, Nothing else."
  },
  "Reasoning": {
    "type": "string",
    "description": "Why you feel the image is from that particular country. Explain what features/objects/items/visuals of the image made you think so. Should be fine-grained and clear, but in just 2-3 sentences."
  },
  "Confidence": {
    "type": "integer",
    "description": "A integer score on a scale of 0 to 100 on how confident you are about the classification."
  },
  "Features": {
    "type": "string",
    "description": "Which of the features of the image were primarily used for determining the country labels, separated by commas, ordered with primary feature in the beginning."
  }
}
```

```
## MCQ
prompt_header="Respond with a label from the choices : What country is this image likely from ?\n"
for l,c in zip(labels,choices):prompt_header+=f"{l}) : {c}\n"
fixed_prompt=prompt_header+"""Respond in the following format :
{
  "Label": "Uppercase Alphabet",
  "Country": "The chosen label's country name exactly as it was",
  "Confidence": "Integer between 0 to 100 in numeric format"
}"""
```

E Other Plots

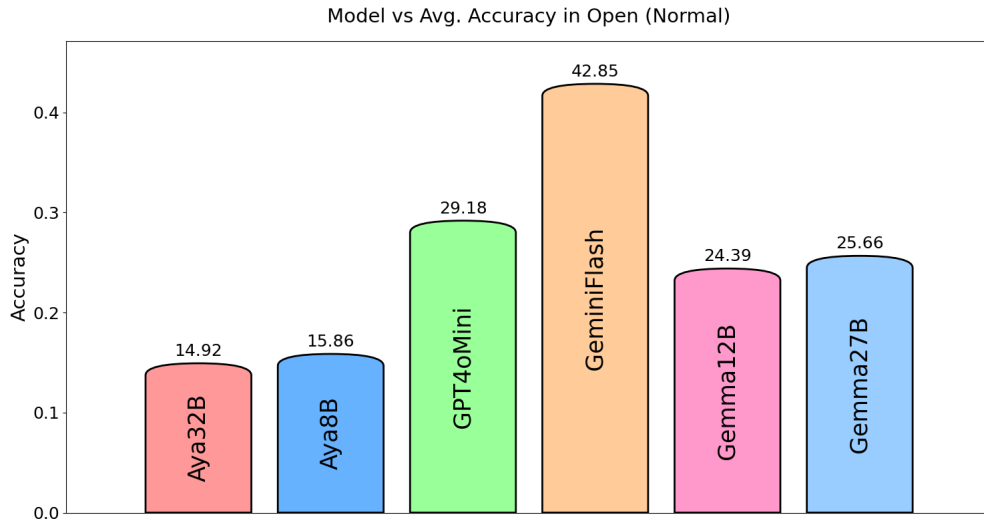


Figure 8: Overall Accuracy : Open Ended (Normal)

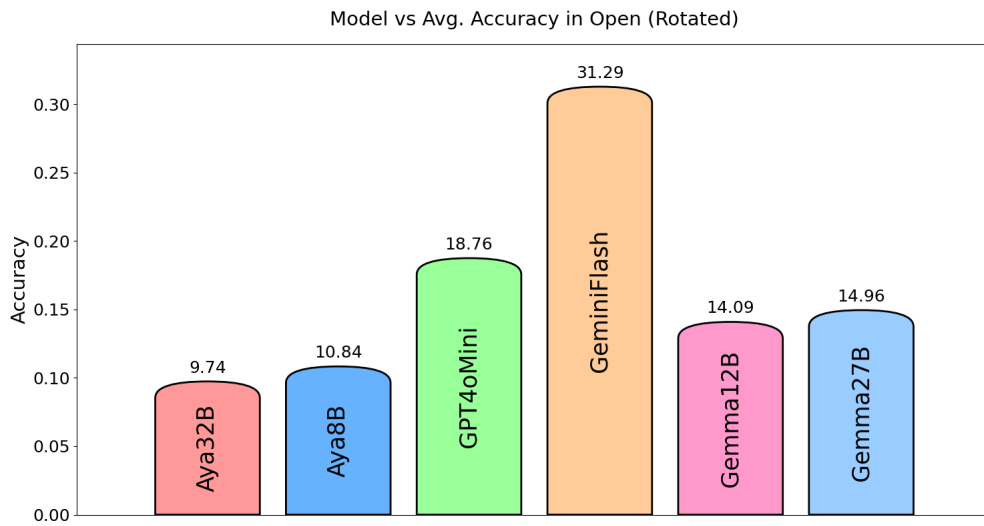


Figure 9: Overall Accuracy : Open Ended (Rotated)

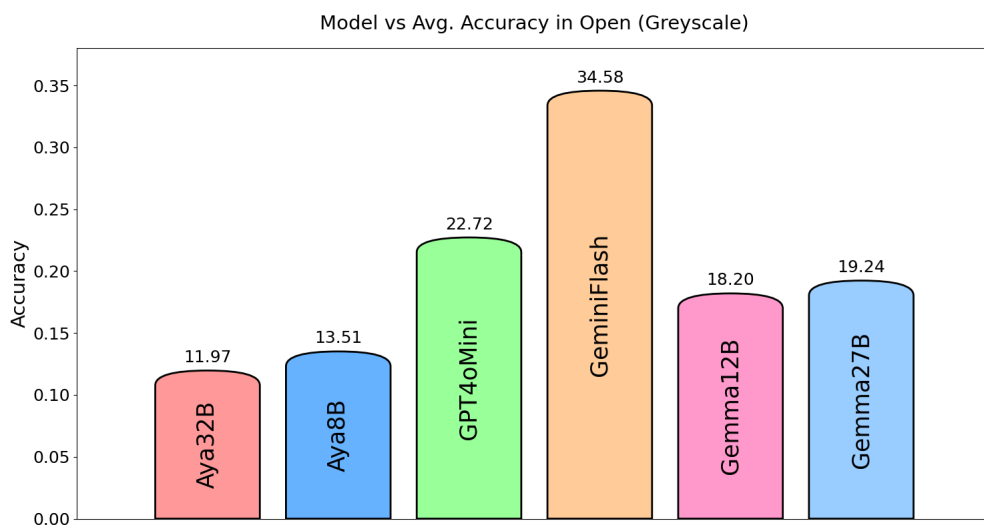


Figure 10: Overall Accuracy : Open Ended (Grayscale)

Category	Description
Appearance (Attire)	Attires of some people from the image, clothes being hanged in the background, etc.
Appearance (People)	Appearance / visual perception of people's ethnicity, presence of any celebrities, etc.
Architecture (Exterior)	Building facades, monuments, bridges, outdoor structures, and any external architectural elements visible in the scene.
Architecture (Interior)	Indoor environments e.g. rooms, corridors, staircases, furniture, and interior design details.
Landscape (Water)	Bodies of water such as oceans, rivers, lakes, waterfalls, ponds, and any aquatic scenery.
Landscape (Air)	Aerial / bird's-eye views, landscapes captured from above, clouds, sky scenes, and horizon vistas.
Landscape (Vegetation)	Forests, grasslands, gardens, crops, shrubs, foliage patterns, plant life, or visible greenery.
Texts/Scripts/Posters	Signs, banners, billboards, labels, handwritten or printed text, posters, and any other written or graphic messaging.
Patterns/Designs	Decorative motifs, surface textures, fabric prints, wallpaper or tile patterns, abstract designs, and repetitive graphical elements.

Table 4: Overview of the image categories used to analyse model performance as a function of the type of image.

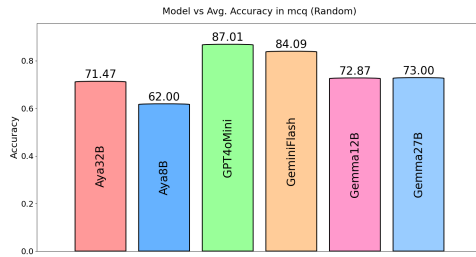


Figure 11: Overall Accuracy : MCQ-Random : Model wise

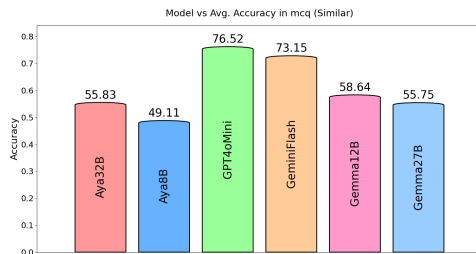


Figure 12: Overall Accuracy : MCQ-Similar : Model wise

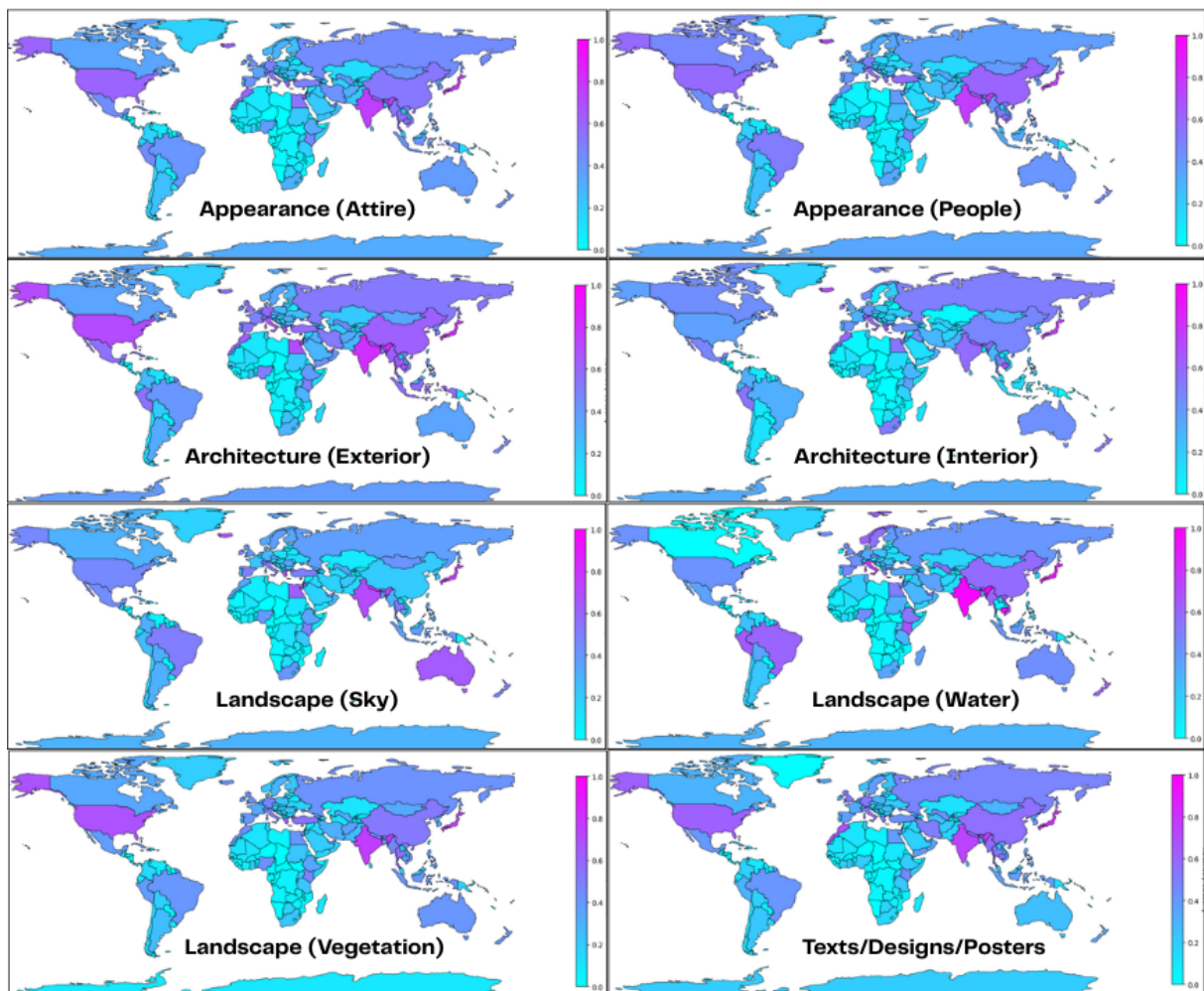


Figure 13: Image Feature categories VS Country wise Accuracy

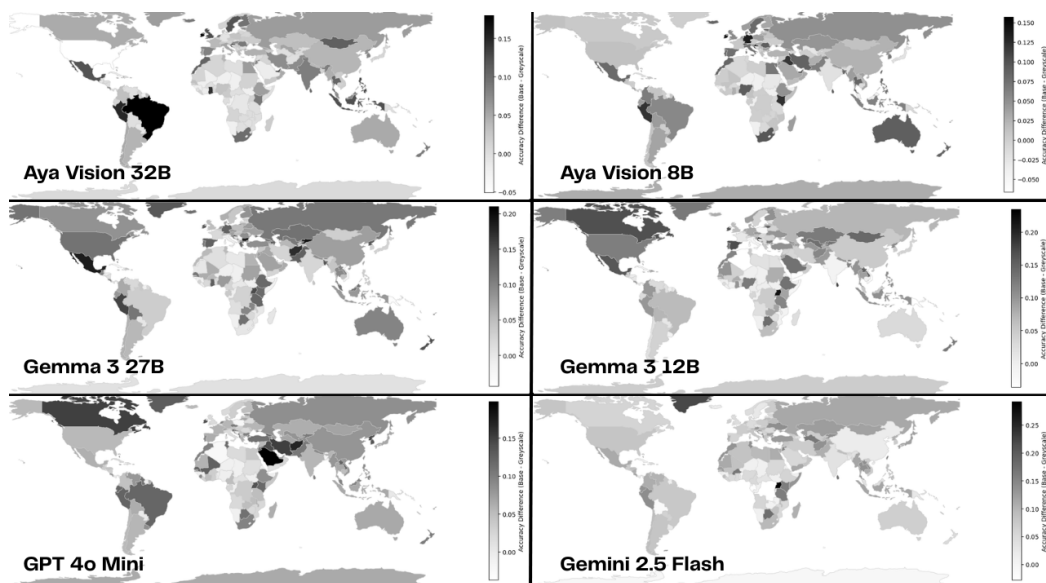


Figure 14: Effect of Gray-scaling VS change in country wise accuracies

Higher Contrast = Larger Drop in accuracy

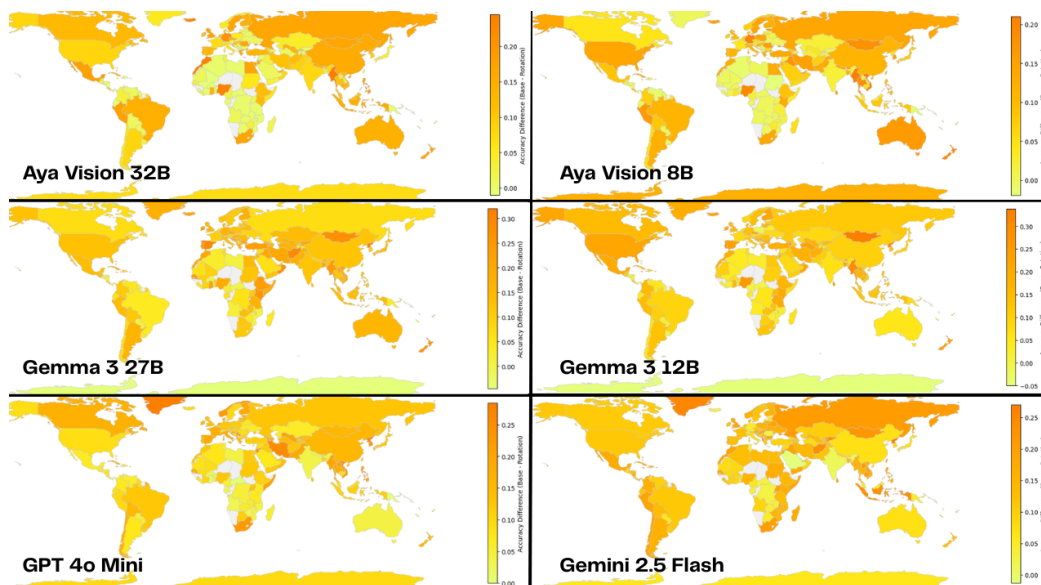


Figure 15: Effect of Rotation VS change in country wise accuracies

Higher Contrast = Larger Drop in accuracy

F Mis-Classification Map : Region-wise

The mis-classifications from one region to countries outside the region can be seen from each region in Figure 20 to Figure 34 respectively.

G Country wise accuracies in each experimental setting

The accuracies obtained over samples of each country through each experimental setup can be seen in Table 5 to Table 9.

Country name	Open-Ended	MCQs with Similar choices	MCQs with Random choices
Afghanistan	41.33	68.90	81.56
Albania	20.00	42.80	67.64
Algeria	10.50	29.73	65.71
Andorra	12.00	59.63	72.41
Angola	4.67	48.07	58.83
Anguilla	2.00	15.27	58.51
Antarctica	34.83	84.80	83.57
Antigua and Barbuda	7.67	31.67	70.64
Argentina	30.67	84.17	71.39
Armenia	42.33	66.23	80.07
Aruba	17.67	55.67	78.96
Australia	44.50	87.90	69.58
Austria	18.83	42.13	80.69
Azerbaijan	20.00	46.83	66.45
Bahamas	24.83	69.47	78.13
Bahrain	21.00	63.00	73.94
Bangladesh	42.50	59.30	87.48
Barbados	17.67	39.50	72.07
Belarus	13.33	45.60	72.98
Belgium	21.00	44.93	72.21
Belize	11.67	59.13	68.49
Benin	7.50	51.47	78.75
Bermuda	20.67	62.63	67.61
Bhutan	59.17	66.03	90.70
Plurinational State of Bolivia	26.33	76.13	78.26
Bonaire, Sint Eustatius and Saba	3.50	36.47	69.24
Bosnia and Herzegovina	23.33	44.43	73.23
Botswana	22.83	82.13	80.00
Brazil	47.67	83.37	74.70
Brunei Darussalam	8.67	21.73	48.78
Bulgaria	25.33	46.47	77.12
Burkina Faso	7.50	60.83	74.72
Cabo Verde	10.17	67.23	55.22
Cambodia	62.83	81.02	92.15
Cameroon	4.67	67.20	70.02
Canada	41.50	69.43	81.16
Cayman Islands	6.67	28.07	68.78
Central African Republic	0.83	16.67	50.21
Chile	20.83	65.90	67.78
China	58.83	78.73	81.48
Colombia	23.83	75.73	69.25
Democratic Republic of Congo	6.83	40.70	56.60
Cook Islands	3.83	22.23	68.28

Table 5: Country wise accuracies through various experimental settings : Part 1/5



Figure 16: Region wise effect of perturbations

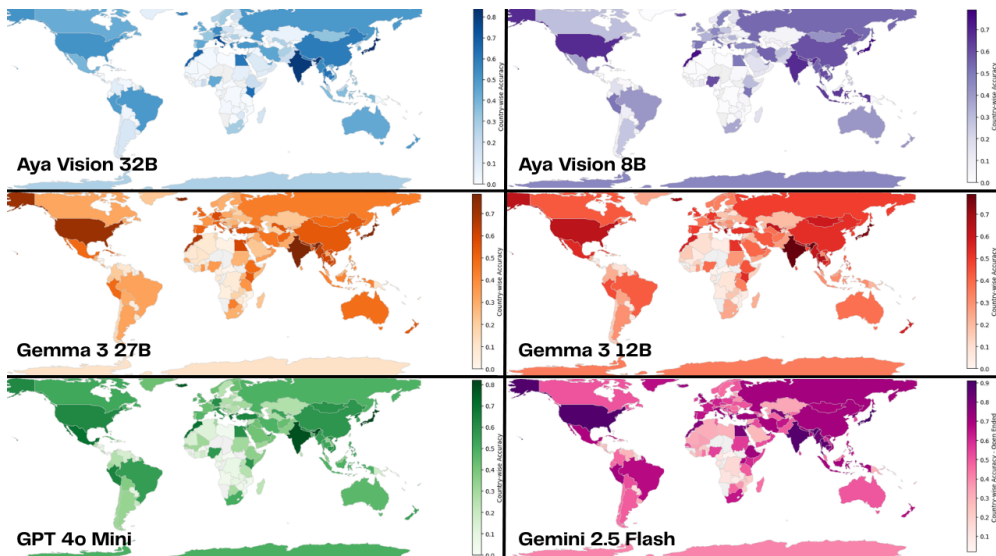


Figure 17: Accuracy over each country's images through open-ended Experiments

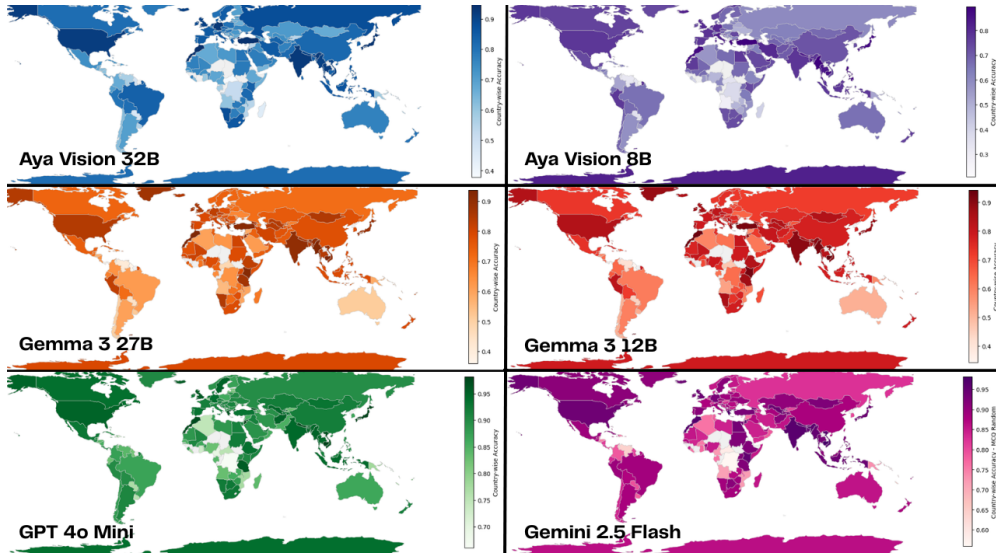


Figure 18: Accuracy over each country's images through MCQ Experiments with random distractors

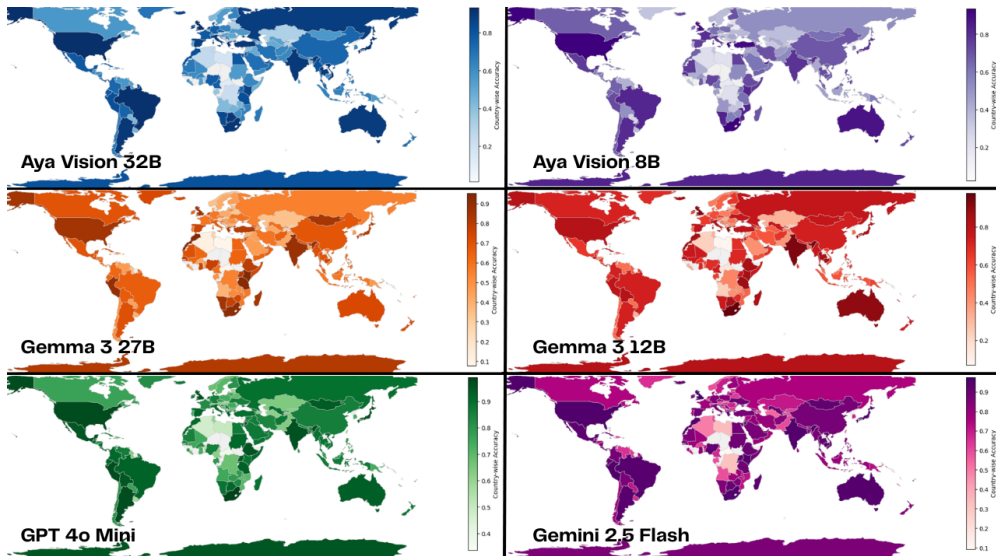


Figure 19: Accuracy over each country's images through MCQ Experiments with similar distractors

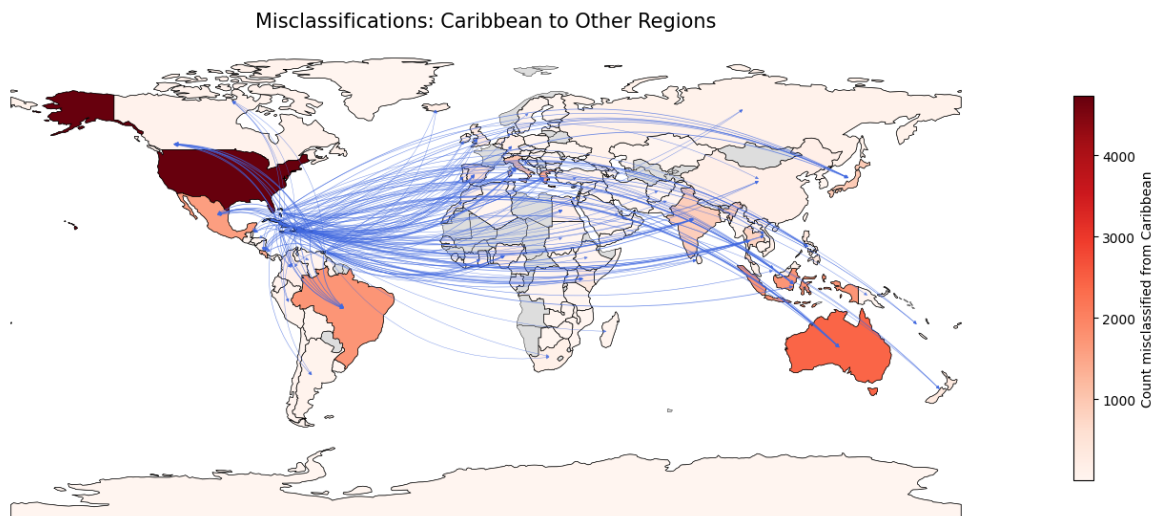


Figure 20: Mis-classification map : Caribbean

Misclassifications: Western Europe to Other Regions

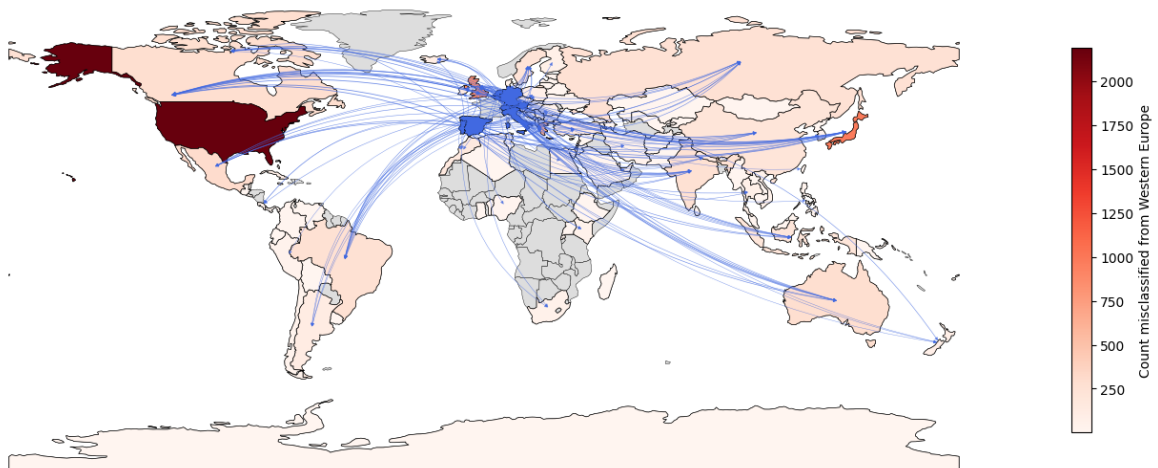


Figure 21: Mis-classification map : Western Europe

Misclassifications: Northern Europe to Other Regions

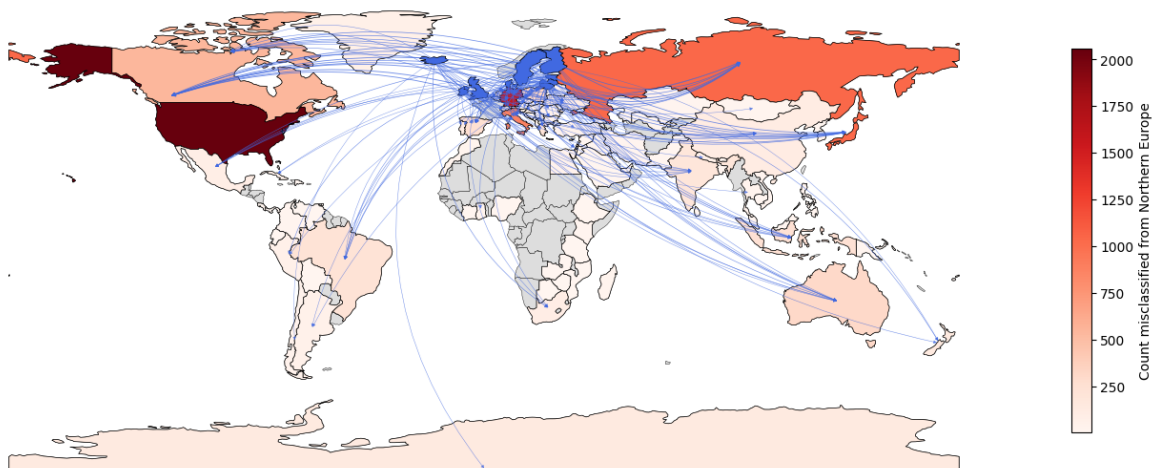


Figure 22: Mis-classification map : North Europe

Misclassifications: Eastern Europe to Other Regions

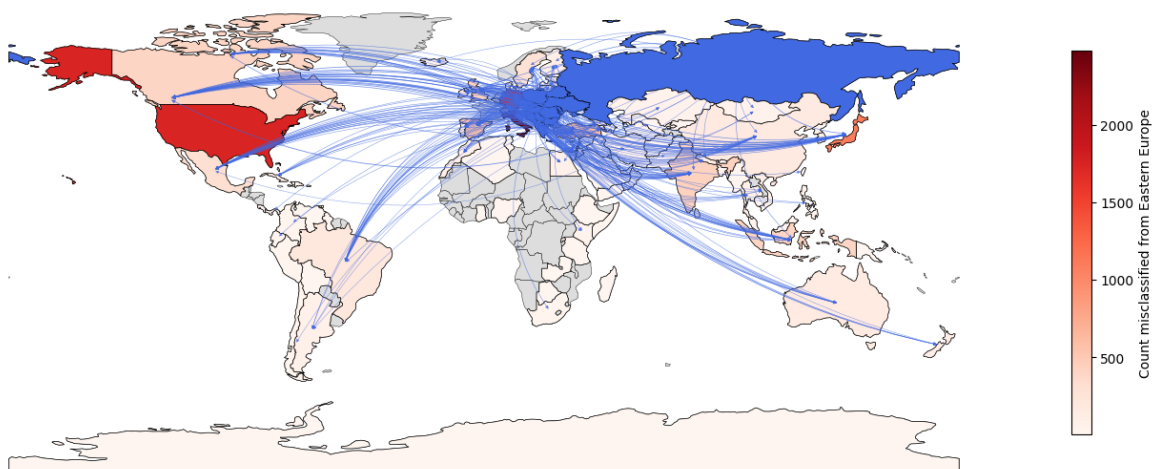


Figure 23: Mis-classification map : Eastern Europe

Misclassifications: East Asia to Other Regions

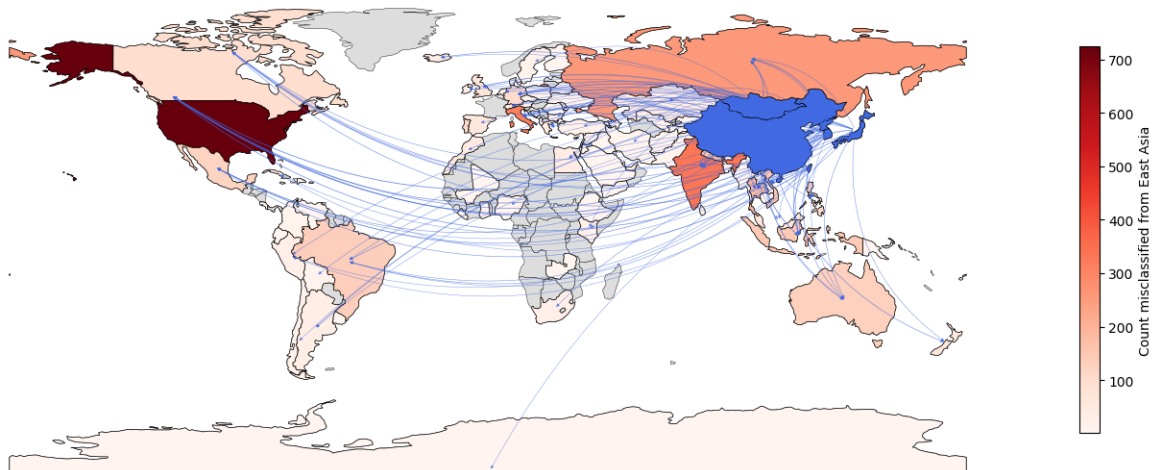


Figure 24: Mis-classification map : East Asia

Misclassifications: Central Asia to Other Regions

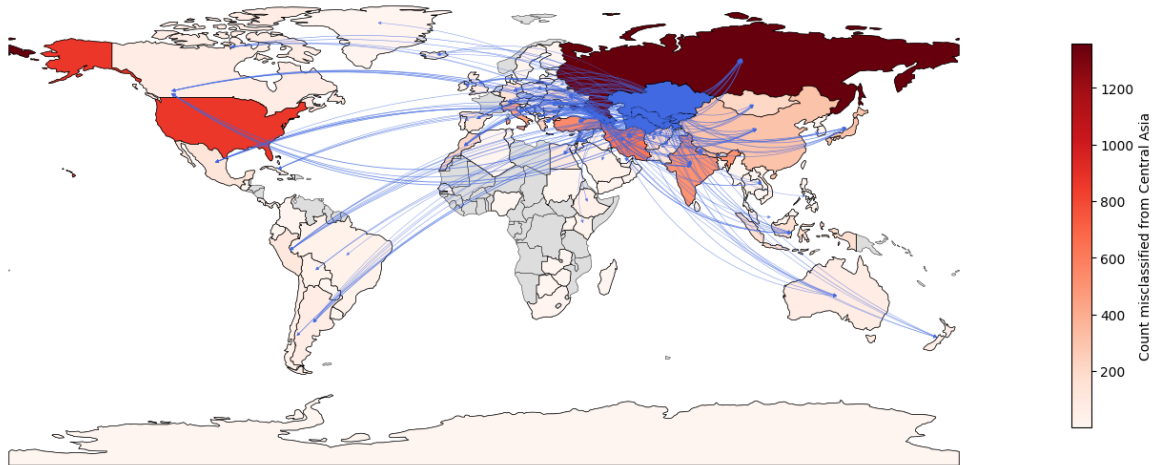


Figure 25: Mis-classification map : Central Asia

Misclassifications: Southeast Asia to Other Regions

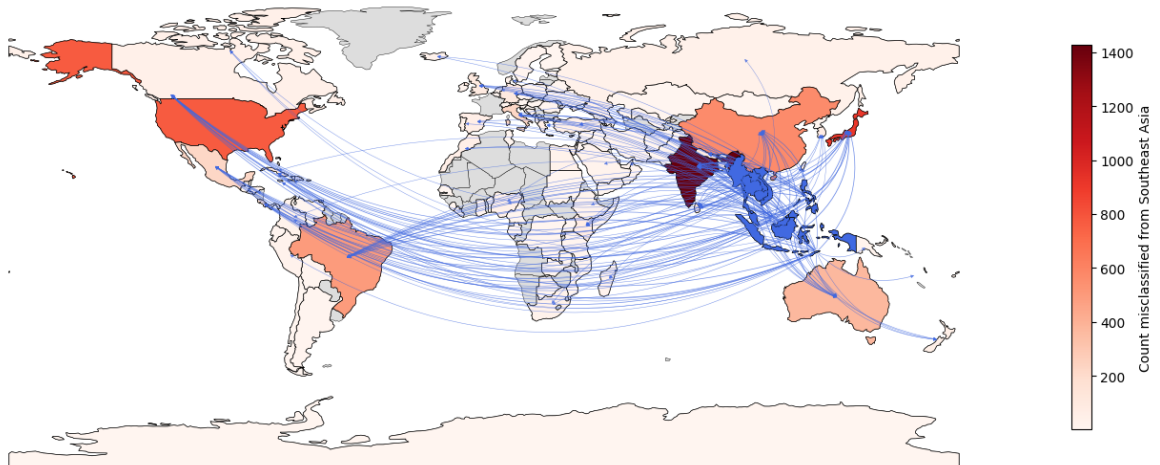


Figure 26: Mis-classification map : South East Asia

Misclassifications: South Asia to Other Regions

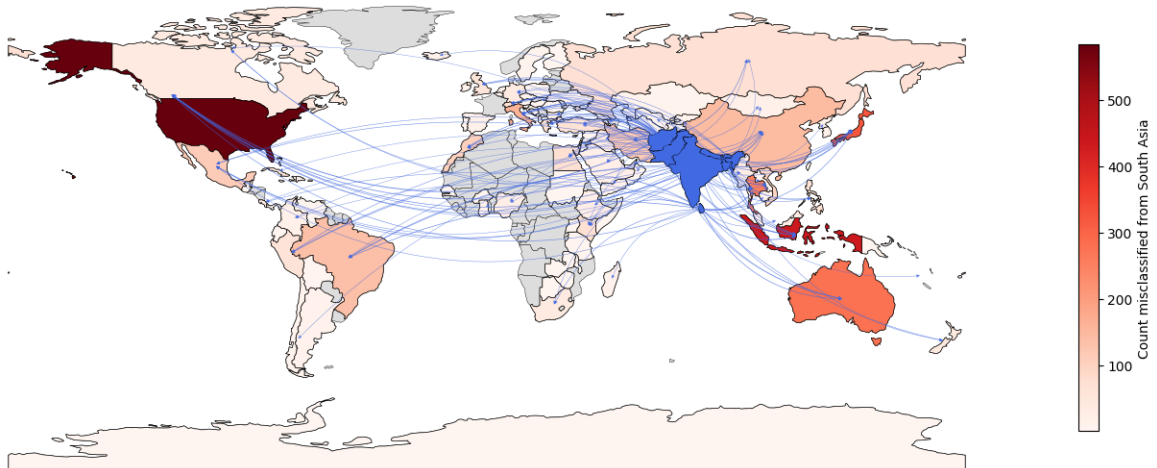


Figure 27: Mis-classification map : South Asia

Misclassifications: Middle East to Other Regions

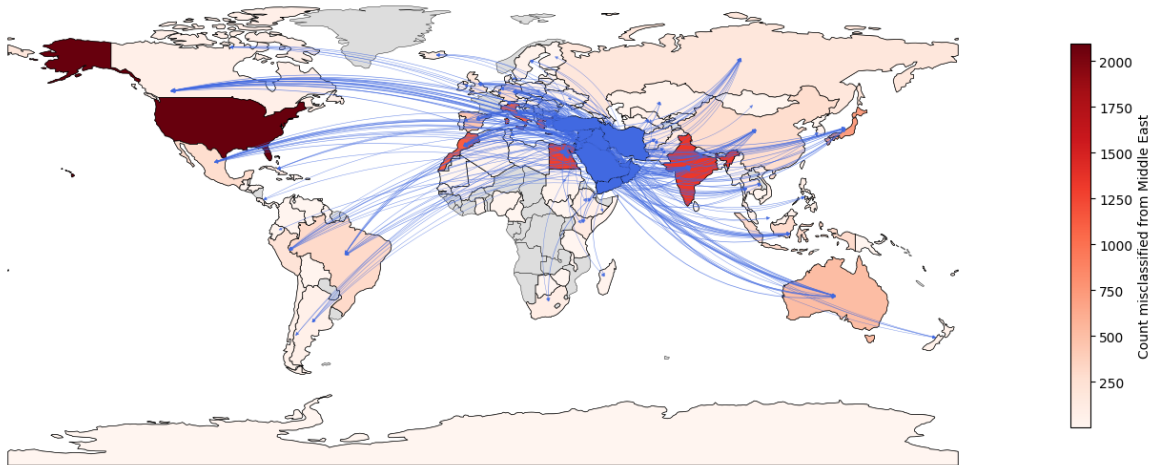


Figure 28: Mis-classification map : Middle East

Misclassifications: Southern Africa to Other Regions

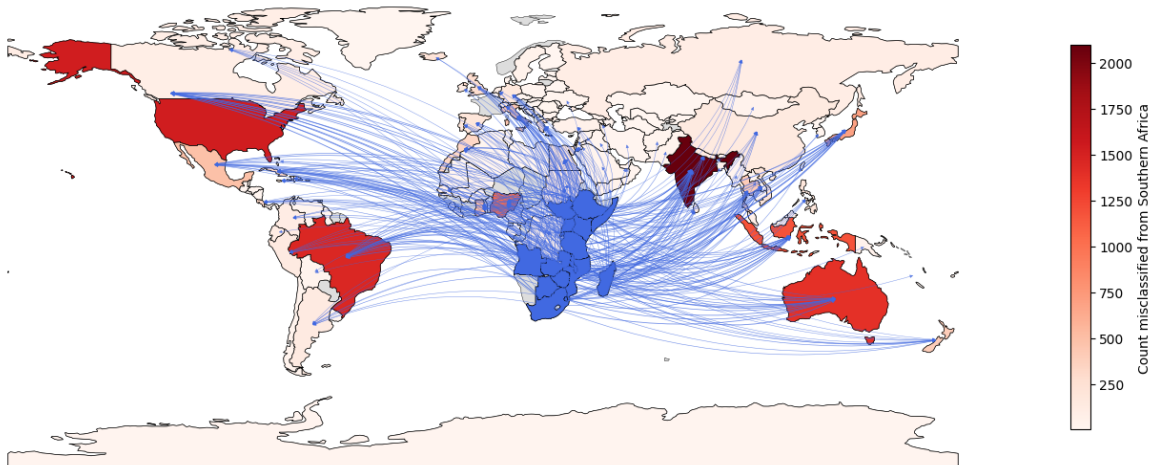


Figure 29: Mis-classification map : Southern Africa

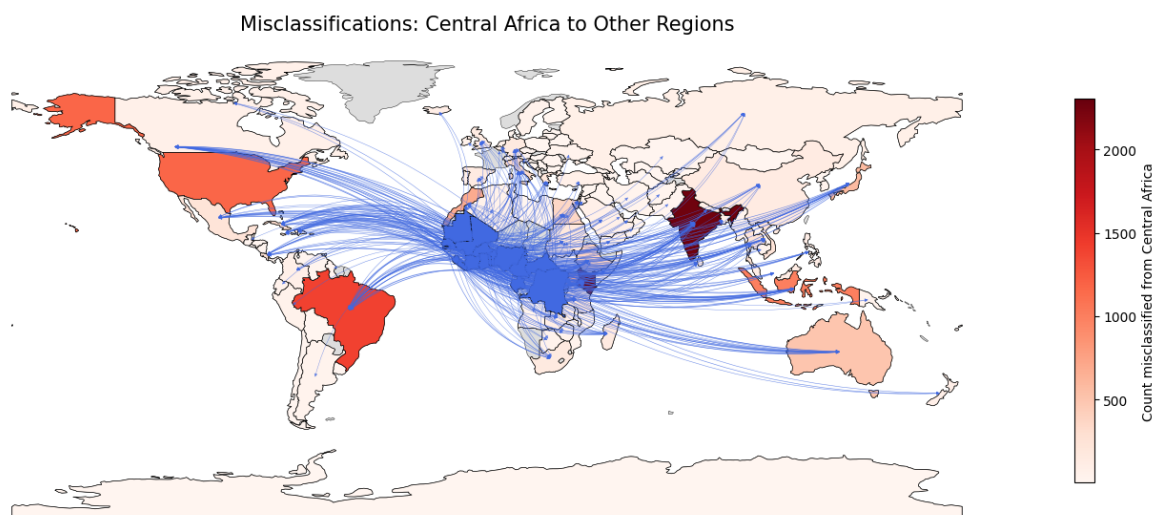


Figure 30: Mis-classification map : Central Africa

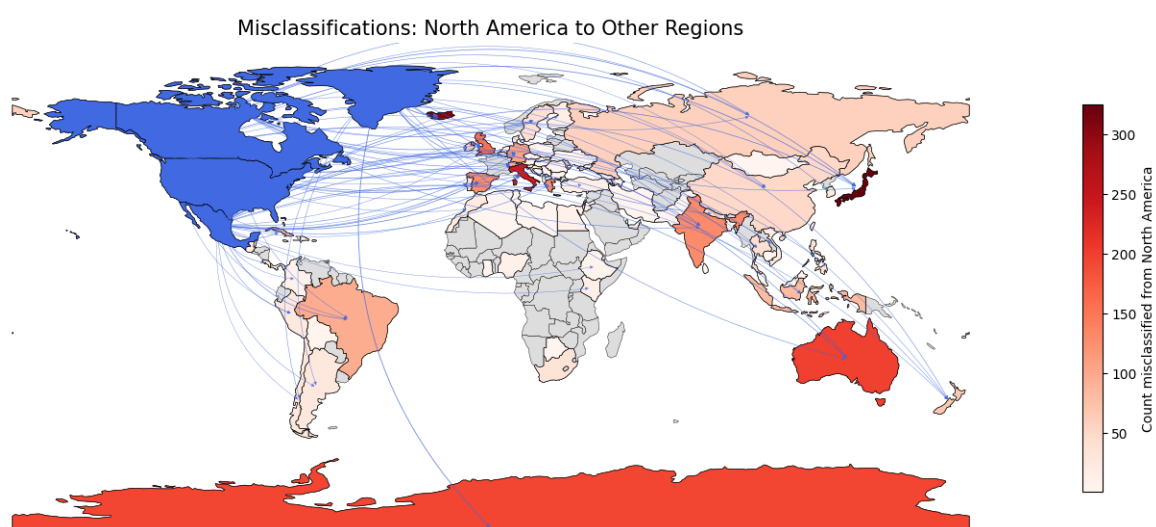


Figure 31: Mis-classification map : North America

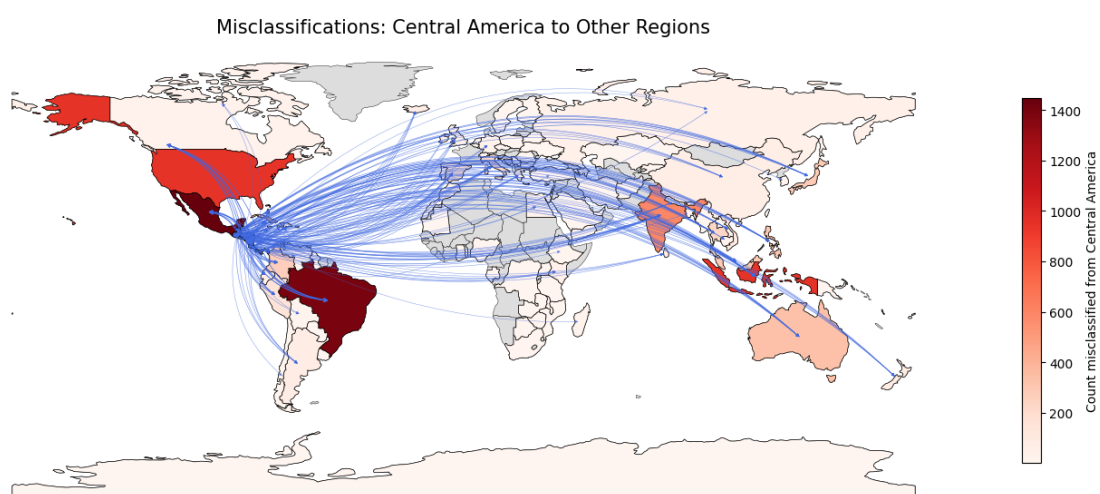


Figure 32: Mis-classification map : Central America

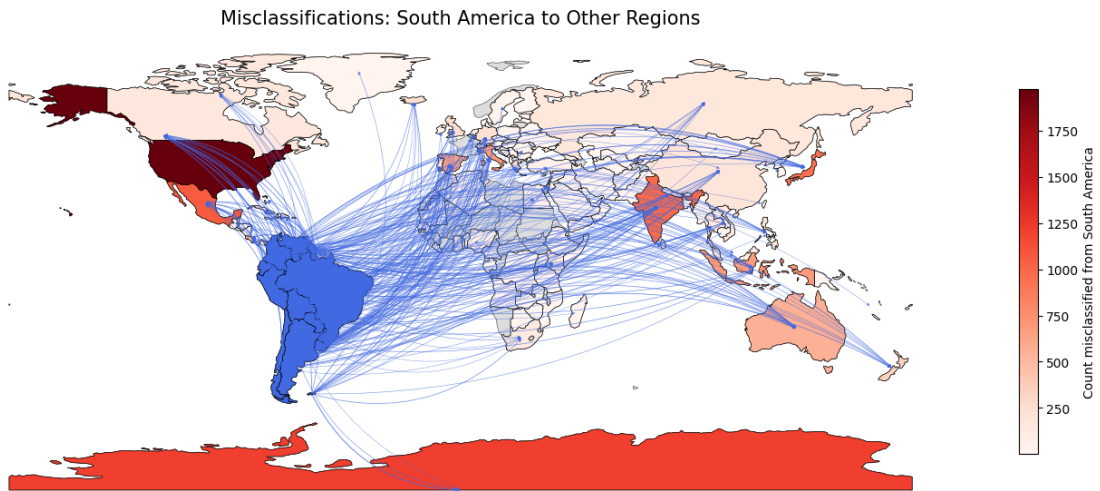


Figure 33: Mis-classification map : South America

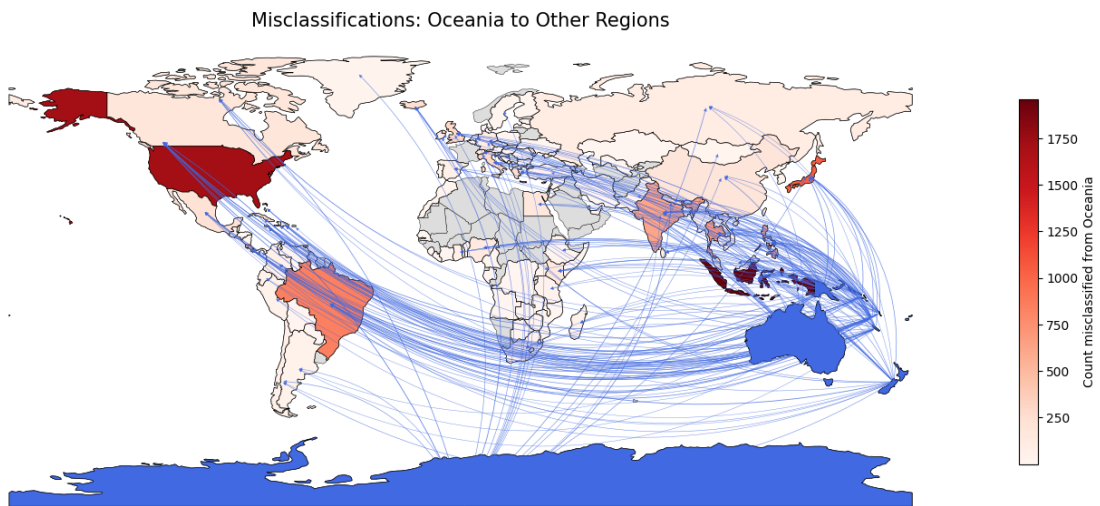


Figure 34: Mis-classification map : Oceania



Figure 35: Examples from ours (1st,4th) as well as other works : GIMMICK (2nd), CVQA (3rd) : The 1st and 4th image have the key features required for classifying the image accurately, occupying a tiny portion of the image making it relatively difficult i.e the flag patch in image 1 ,and name of mountain in image 4's signboard. While, in Image 2 and image 3 , the key features i.e the text on attire or the (car, city name signboard, multilingual texts on left) make the samples relatively easier to classify

Country name	Open-Ended	MCQs with Similar choices	MCQs with Random choices
Costa Rica	26.00	73.23	72.16
Croatia	47.83	72.83	83.92
Cuba	47.50	76.83	77.92
Curaçao	20.83	61.07	80.96
Cyprus	13.67	59.33	69.19
Czechia	40.50	66.07	83.90
Côte d'Ivoire	13.33	60.00	71.47
Denmark	32.50	66.93	78.54
Dominica	15.17	61.17	67.04
Dominican Republic	15.00	56.37	70.43
Ecuador	21.50	76.10	73.05
Egypt	60.50	77.07	83.84
El Salvador	4.83	65.93	63.40
Estonia	21.83	43.90	70.30
Eswatini	0.50	28.70	53.07
Ethiopia	41.00	80.93	80.24
Falkland Islands	8.83	92.13	90.35
Faroe Islands	30.33	71.30	90.66
Fiji	22.83	64.10	76.93
Finland	32.33	67.80	76.31
France	40.83	73.77	83.70
French Guiana	3.00	64.73	53.93
French Polynesia	24.67	81.90	83.97
Gabon	5.33	56.00	66.67
Gambia	3.33	41.40	53.16
Georgia	32.00	71.40	83.37
Germany	54.83	71.30	87.54
Ghana	26.33	70.53	67.80
Gibraltar	19.00	62.27	79.48
Greece	66.67	91.00	91.06
Greenland	27.00	65.43	84.90
Grenada	3.50	37.77	63.75
Guadeloupe	1.50	48.80	71.09
Guam	11.33	70.43	55.57
Guatemala	19.67	75.50	74.38
Guernsey	1.83	66.50	81.14
Guyana	8.83	52.33	52.66
Haiti	27.83	71.23	65.98
Vatican City State	8.67	43.77	74.31
Honduras	4.67	64.13	66.98
Hong Kong	22.67	65.23	86.70
Hungary	24.83	49.00	78.44
Iceland	69.00	82.27	89.04

Table 6: Region wise accuracies through various experimental settings : Part 2/5

Country name	Open-Ended	MCQs with Similar choices	MCQs with Random choices
India	78.33	90.03	90.10
Indonesia	48.83	67.76	84.97
Iran	50.83	70.40	83.27
Iraq	28.67	60.60	76.84
Ireland	48.33	74.57	87.63
Isle of Man	6.17	52.03	77.91
Israel	35.67	76.33	73.99
Italy	60.00	82.30	85.40
Jamaica	28.17	60.20	70.58
Japan	81.17	88.92	91.75
Jersey	3.67	50.37	71.69
Jordan	44.00	79.03	89.04
Kazakhstan	18.33	44.73	77.73
Kenya	56.00	88.57	88.15
North Korea	47.33	25.64	81.26
South Korea	47.83	67.23	79.90
Kuwait	12.83	52.30	68.70
Kyrgyzstan	20.17	37.30	69.48
Laos	26.50	38.53	80.25
Latvia	17.00	41.63	72.47
Lebanon	27.00	73.63	78.09
Liberia	9.33	50.97	65.37
Libya	6.67	22.87	73.10
Liechtenstein	6.17	34.03	72.29
Lithuania	24.00	54.43	74.40
Luxembourg	13.33	21.90	62.29
Macao	17.00	66.42	85.38
Madagascar	24.17	81.20	65.40
Malawi	8.33	54.80	66.39
Malaysia	28.33	73.28	83.65
Maldives	39.33	80.20	82.08
Mali	13.83	65.43	80.11
Malta	47.67	79.57	90.95
Martinique	4.33	53.60	72.85
Mauritania	12.00	76.77	80.28
Mauritius	38.33	92.00	79.52
Mexico	53.17	79.77	79.69
Moldova	7.67	35.23	63.57
Monaco	30.17	54.83	69.69
Mongolia	50.83	82.41	81.39
Montenegro	22.17	44.37	81.00
Morocco	67.83	85.40	93.75
Mozambique	5.17	66.57	63.78
Myanmar	61.50	76.56	92.62

Table 7: Region wise accuracies through various experimental settings : Part 3/5

Country name	Open-Ended	MCQs with Similar choices	MCQs with Random choices
Namibia	0.00	83.40	85.35
Nepal	65.00	72.53	89.72
Netherlands	46.00	74.63	86.86
New Caledonia	7.50	55.03	64.98
New Zealand	53.83	76.40	82.58
Nicaragua	6.83	69.87	69.64
Nigeria	47.33	79.13	73.78
North Macedonia	10.17	44.27	74.44
Norway	32.50	48.17	79.45
Oman	31.67	71.40	77.59
Pakistan	30.33	53.57	79.32
Palau	15.83	71.23	71.97
Palestine, State of	9.00	73.53	83.59
Panama	4.33	80.17	60.86
Papua New Guinea	13.50	61.87	63.38
Paraguay	6.17	52.23	54.29
Peru	54.83	85.73	83.61
Philippines	43.67	74.82	85.94
Poland	28.83	62.00	79.17
Portugal	43.50	58.60	84.39
Puerto Rico	16.67	68.97	72.52
Qatar	19.50	56.63	66.04
Romania	31.50	56.43	79.02
Russian Federation	52.67	73.13	77.18
Rwanda	29.50	71.73	73.72
Réunion	5.33	90.87	69.21
Saint Helena, Ascension and Tristan da Cunha	3.33	71.40	57.44
Saint Kitts and Nevis	14.17	41.23	64.61
Saint Lucia	16.83	61.40	79.33
Saint Martin (French)	4.00	45.43	69.48
Samoa	23.33	68.43	71.19
San Marino	10.17	35.00	54.01
Saudi Arabia	26.00	65.53	74.69
Senegal	21.83	78.73	78.20
Serbia	24.33	58.70	79.14
Seychelles	26.33	92.87	76.83
Sierra Leone	8.83	56.53	75.23
Singapore	51.33	74.91	80.15
Saint Martin (Dutch)	7.17	50.77	75.14
Slovakia	12.33	32.33	67.41
Slovenia	24.00	53.40	75.09
Solomon Islands	3.33	22.53	69.22
Somalia	24.67	75.30	78.46

Table 8: Region wise accuracies through various experimental settings : Part 4/5

Country name	Open-Ended	MCQs with Similar choices	MCQs with Random choices
South Africa	38.50	94.43	82.91
South Georgia and the South Sandwich Islands	7.17	80.70	77.99
South Sudan	25.83	65.83	82.31
Spain	51.00	83.13	84.71
Sri Lanka	37.00	61.40	82.72
Sudan	25.33	70.63	81.25
Svalbard and Jan Mayen	0.00	74.13	89.45
Sweden	35.50	54.63	81.22
Switzerland	42.17	62.53	76.40
Syrian Arab Republic	13.00	51.63	64.82
Taiwan, Province of China	23.00	51.01	80.16
Tajikistan	10.83	44.43	81.04
Tanzania, United Republic of	24.83	84.37	84.89
Thailand	64.17	84.49	89.08
Timor-Leste	7.83	41.77	69.67
Togo	2.33	31.67	65.98
Tonga	1.33	19.60	44.73
Trinidad and Tobago	8.00	56.23	53.62
Tunisia	20.33	40.00	75.53
Turkmenistan	22.67	48.73	82.83
Türkiye	56.33	86.10	92.24
Uganda	26.83	79.90	80.27
Ukraine	22.83	67.63	72.82
United Arab Emirates	53.00	85.30	85.30
United Kingdom	50.17	92.17	89.05
United States	67.17	91.03	87.76
Uruguay	14.17	46.33	61.10
Uzbekistan	47.17	68.63	83.07
Vanuatu	5.50	18.00	57.04
Venezuela, Bolivarian Republic of	11.17	57.63	53.41
Viet Nam	55.50	78.74	89.77
Virgin Islands, British	6.83	38.00	79.60
Virgin Islands, U.S.	9.67	46.73	81.72
Kosovo	6.50	28.70	65.53
Yemen	27.17	69.80	76.46
Zambia	9.50	54.80	73.29
Zimbabwe	11.67	71.03	76.05
Åland Islands	0.17	29.00	62.02
Overall	25.14	61.92	75.06

Table 9: Region wise accuracies through various experimental settings : Part 5/5



Country : Nepal

Confidence : 85

Features : ['prayer flags', 'architectural style', 'ornate windows']

Reasoning : The presence of colorful prayer flags is a strong indicator of Nepalese culture, commonly found in Buddhist and Hindu contexts. Additionally, the architectural style of the buildings in the background reflects traditional Nepali design, which further supports the classification.

Figure 36: A sample from our dataset and its corresponding response (GPT-4o-Mini)