048

049

050

051

052

053

054

055

056

057

058

059

CultureVLM: Characterizing and Improving Cultural Understanding of Vision-Language Models for over 100 Countries

Anonymous CVPR submission

Paper ID 00006

Abstract

001 Vision-language models (VLMs) have advanced human-AI 002 interaction but struggle with cultural understanding, often misinterpreting symbols, gestures, and artifacts due to bi-003 004 ases in predominantly Western-centric training data. In this paper, we construct CultureVerse, a large-scale mul-005 006 timodal benchmark covering 19,682 cultural concepts, 188 countries/regions, 15 cultural topics, and 3 question types, 007 008 with the aim of characterizing and improving VLMs' multicultural understanding capabilities. Then, we propose 009 010 CultureVLM, a series of VLMs fine-tuned on our dataset to achieve significant performance improvement in cultural 011 understanding. Our evaluation of 16 models reveals sig-012 nificant disparities, with a stronger performance in West-013 ern concepts and weaker results in African and Asian con-014 015 texts. Fine-tuning on our CultureVerse enhances cultural 016 perception, demonstrating cross-cultural, cross-continent, and cross-dataset generalization without sacrificing perfor-017 mance on models' general VLM benchmarks. We further 018 present insights on cultural generalization and forgetting. 019 020 We hope that this work could lay the foundation for more equitable and culturally aware multimodal AI systems. 021

1. Introduction

Vision-language models (VLMs) have achieved great performance in various tasks, such as visual question answering and captioning [6, 20, 33, 46, 55, 57]. Meanwhile, one
of the most vital aspects of human experience, cultural understanding, which encompasses language, cultural values,
social norms, culinary practices, and artistic expressions–
remains a challenging area for these models [3, 59].

Challenges. Cultural understanding is essential for AI systems intended for global deployment, as it enables them
to interact appropriately and sensitively with users of diverse cultural, ethnic, and social backgrounds. However,
current VLMs often struggle to grasp the deeper cultural
meanings embedded in symbols and artifacts. For instance,

a VLM may identify an eagle as merely a bird, overlooking 036 its symbolic significance as a national emblem representing 037 the spirit and identity of the United States. Similarly, the 038 lotus flower is not only a plant, but a profound symbol of 039 purity and spiritual enlightenment in Indian culture. Ges-040 tures present an even more complex challenge: the "OK" 041 hand gesture, which conveys a positive meaning in North 042 American countries, is interpreted as offensive in countries 043 such as Brazil and Turkey [35]. Misinterpretations of cul-044 turally significant symbols can lead to misunderstandings 045 and even cause offense. 046

These challenges partially stem from inherent biases and limitations in VLMs' training data: 1) Skewed Domain Coverage. Pre-training images and texts predominantly feature generic daily scenes or natural settings, often lacking coverage of culturally specific artifacts, traditions, beliefs, and historical sites. Models may fail to interpret culturally significant symbols, particularly those from underrepresented regions. 2) English-centric Data and Western Bias. The texts for pre-training VLMs is primarily sourced from English content [21, 39], which predominantly represents high-resource cultures, resulting in a Western Bias [15, 61]. This limits the models' understanding of diverse cultures, especially those in the global south [13].

Addressing these challenges is crucial for developing 060 culturally aware AI systems that can engage effectively and 061 respectfully with global users. Although there have been 062 efforts to build culturally aware LLMs through data collec-063 tion [13, 27, 52], improving VLMs for cultural understand-064 ing remains in its infancy. VLMs require multimodal input-065 both texts and images-making the collection or generation 066 of culturally rich training data even more challenging, es-067 pecially for low-resource cultures. From the benchmark-068 ing perspective, existing datasets for VLMs [41, 50] usually 069 rely on human annotators for data curation. These datasets 070 are limited in scale, often lack sufficient regional and na-071 tional representation, and may not capture deep cultural rel-072 evance. From the *modeling* perspective, there is a lack of 073 work aimed at building culturally aware VLMs, a signifi-074 cant obstacle to AI equity for underrepresented cultures. 075

CVPR 2025 Submission #00006. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.



Figure 1. Our pipeline to build CultureVerse and CultureVLM.

076 This Work. We take the first step toward advancing cultural understanding in VLMs through both comprehen-077 sive benchmarking and targeted model improvements. To 078 achieve this, we construct CultureVerse, a large-scale mul-079 timodal dataset to evaluate and enhance the multicultural 080 capabilities of VLMs. Our flexible pipeline (Figure 1) can 081 easily integrate additional languages and cultures, ensuring 082 adaptability and inclusiveness. Our work lays a solid foun-083 dation for developing more equitable AI systems that ad-084 dress the needs of developing nations, ethnic minorities, and 085 086 underrepresented cultures. The key findings are as follows.

- Disparity in Cultural Understanding: *All* VLMs show highly consistent regional disparities in cultural understanding, with the highest cultural understanding for the Americas, followed by Europe and Oceania, and the weakest understanding for Asia and Africa.
- Training for Enhanced Cultural Perception: Finetuning effectively enhances the cultural perception of VLMs, narrowing the gaps in cultural understanding across different regions and categories without significantly compromising the model's general capabilities.
- Model and Data Scale Enhance Cultural Understanding: Cultural understanding is generally positively correlated with model size, though not absolutely, as demonstrated by the Llama 3.2-11B model achieving performance comparable to that of Qwen 2-72B. Regarding fine-tuning, larger training datasets lead to more significant improvements, but gains slow as data grows.
- Generalization across Cultures, Concepts, Continent, and Datasets: Due to the inherent correlations between cultures of different regions and types, fine-tuning for cultural understanding exhibits reasonable generalization across different cultures, concepts, continents, and even datasets, showing great potential to improve cultural understanding via generalization research.
 - **Contributions.** Our contributions are as follows:

111

- Large-scale Dataset. We present CultureVerse, a massive scale benchmark consisting of 19, 682 cultural concepts and 228, 053 samples, covering 3 tasks, 188 countries or regions, and 15 cultural topics. The test set includes 11, 085 widely recognized concepts and their corresponding 31, 382 samples.
- **Comprehensive Evaluation.** We evaluate a wide range

of cultural concepts across 16 open-source and proprietary models of varying scales. 120

 Improvement of Multimodal Cultural Understanding. We present CultureVLM, which includes a flexible and cost-effective data collection and construction process and a series of VLMs fine-tuned on our dataset. Experimental results demonstrate that our models enhance cultural understanding while maintaining general capabilities and exhibiting a degree of generalization abilities.
 121 122 123 124 125 126 127

2. CultureVerse: A Scalable Benchmark for VLM Cultural Understanding 128

Collecting reliable large-scale culture datasets presents two 130 key challenges: Diversity and Scalability. Achieving com-131 prehensive coverage is especially difficult for culturally di-132 verse topics, particularly for underrepresented groups in the 133 Global South. The construction and annotation of such 134 datasets would typically require substantial human exper-135 tise from various countries and ethnic groups, resulting in 136 poor scalability and high costs. Existing cultural bench-137 marks for VLMs usually lack adequate representation of 138 diverse regions and communities, and often reflect a bias 139 towards dominant cultures [31, 44, 50]. 140

To overcome these limitations, we introduce a scal-141 able data collection pipeline that integrates automated web 142 crawling for scalability and diversity with expert human an-143 notation for *reliability*. As shown in Figure 1, our pipeline 144 consists of three stages: tangible cultural concept collection 145 (2.1), question-answer generation (2.2), and quality assur-146 ance (B.1). This hybrid approach ensures that our dataset 147 captures a wide spectrum of cultural contexts while main-148 taining high standards of data quality and relevance. At the 149 same time, this pipeline is more scalable (B.2) compared to 150 existing methods, regardless of concepts, language, country, 151 region, or image. More details and analyses are in B.3. 152

2.1. Tangible Cultural Concept Collection

To construct a comprehensive set of cultural concepts, a154common approach is to employ a bottom-up strategy that155retrieves specialized knowledge from open web documents.156For example, Fung et al. [16] begin with an initial set157of cultural topics (e.g., education, marriage customs, and158

200

201



Figure 2. Overview of CultureVerse. In total, there are over 220k instances and 19k cultural concepts for training and evaluation, respectively, composed of 3 different types of questions from 188 countries.

holiday traditions), collect relevant Wikipedia documents,
and expand their scope through linked connections. However, many resulting documents primarily describe general,
abstract, or high-level concepts, such as Renaissance
Art or Mediterranean cuisine, which often lack
specific, unambiguous visual representations.

Concept Construction. To overcome this issue, we adopt 165 166 a top-down approach, starting with 15 predefined cultural 167 categories of tangible cultural concepts such as food, festivals, landmarks, and performing arts, as shown in Table 2. 168 These categories are chosen to capture culturally distinc-169 170 tive and visually recognizable elements suitable for image 171 retrieval and analysis. We then use GPT-40 to process all 172 relevant Wikipedia documents, extracting conceptual enti-173 ties that align with the 15 predefined categories. To ensure the quality and specificity of the extracted entities, we im-174 plement a 3-step filtering process on the extracted concep-175 tual entities: 1) Entity Consolidation. We unified duplicate 176 177 entities, merging those that are identical or differ solely by case, and eliminated entities with formatting issues or ir-178 regularities; 2) Frequency-based Thresholding. We retained 179 only entities that appeared at least twice across the docu-180 ments from a given country, ensuring that the concepts are 181 well-recognized within their cultural context; 3) Entity Re-182 183 finement. We filtered out overly abstract or generic entities such as Imperial Cuisine and those lacking distinct 184 regional specificity such as Steak using additional judg-185 ment by GPT-40. Through this refined process, we curated a 186 187 collection of over 19,682 cultural concepts from 188 coun-188 tries, as shown in Table 10. Our pipeline ensures that the 189 selected concepts are diverse and recognizable, relevant for evaluating VLMs in understanding global cultural diversity. 190 Image Retrieval. Using these concepts and their corre-191 sponding countries, we scrape images from Google Images 192 for each cultural concept¹, obtaining five images for each 193 concept. The first image was reserved for the test set and for 194 human quality assessment, while the remaining four images 195 were used for the training set. Images larger than 10MB 196 were compressed to ensure compatibility with typical input 197 requirements of VLMs. 198

2.2. Question-Answer Generation

We designed three levels of VQA tasks to assess and improve the multicultural knowledge of VLMs:

Image Recognition Questions evaluate models' ability to 202 identify cultural concepts in images. Accurate identifica-203 tion of such concepts is fundamental to retrieving relevant 204 cultural knowledge. Given an image and a cultural concept, 205 models answer questions like "What dish is in the image?" 206 Cultural Knowledge Questions further evaluate model's 207 deeper understanding of the cultural background associated 208 with the concepts. For each concept, we generated compre-209 hensive descriptions, including aspects like location, char-210 acteristics, history, and cultural significance. Subsequently, 211 we instruct GPT-40 to formulate a question based on the in-212 troduction and the image to probe this cultural knowledge 213 without directly naming the concept in the question. These 214 questions require the model to identify cultural concepts 215 and apply various levels of reasoning, drawing on relevant 216 cultural knowledge to provide accurate answers. 217

¹All images are only used for research purpose.

238

239

240

241

242

243

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

Scene Understanding Ouestions are designed to assess the 218 model's ability to interpret, interact, and respond within cul-219 220 turally specific contexts, rather than simply recalling factual information as in the previous two categories. We 221 222 curate scenarios with cultural elements or characteristics depicted in the images, providing context cues that chal-223 lenge the VLMs to make contextually appropriate choices. 224 Using the detailed introductions of the previous steps, we 225 226 prompt GPT-40 to generate scenario-based understanding questions. These questions require the model to not only 227 228 recognize cultural concepts but also apply contextual rea-229 soning based on the associated cultural knowledge.

3. Experiments

3.1. Main Results on CultureVerse

We evaluate our benchmark on 14 open-source and 2 proprietary VLMs. The experimental setup is in C. As shown in Figure 3, our main findings are as follows:

- Task Characteristics: Cultural Scene Understanding
 Outperforms Recognition. (D.1)
 - **Regional Disparities**: Better Performance in Western Cultures. (D.2)
 - Category Characteristics: Weak Understanding of History and Landmarks. (D.3)
 - **Performance Variability from Model Level**: Although larger models tend to demonstrate better performance, size alone is not the determining factor. (D.4)
- Direct Answer v.s. Stepwise Reasoning: Stepwise reasoning does not improve cultural recognition and, in most cases, significantly impairs performance. (D.5)
- 247 More detailed results are in Appendix D.

248 3.2. Training CultureVLM

We fine-tuned three models: LLaVA-1.5, Phi-3-Vision, and 249 LLaMa-3.2-Vision, with the results presented in Figure 4b. 250 251 Cultural knowledge, in contrast to reasoning tasks such as 252 mathematics and coding, is relatively easier to enhance as a 253 form of memory-based perception. Consequently, all four models achieved consistent and substantial improvements, 254 reaching performance levels comparable to those of closed-255 256 source models. We also performed ablation studies to ana-257 lyze the impact of fine-tuning data size (E.1) and decoding 258 temperatures (E.2).

259 3.3. Generalization and Robustness

To evaluate the generalizability of VLMs in multicultural
contexts, we partitioned the training data by continents
(Americas, Asia, Europe, Africa, Oceania) and fine-tuned
LLaVA-1.5-7B [33] separately on each subset. Table 5
in the appendix illustrates the performance of each model
trained on specific continental data and tested in all regions,
and Figure 6 shows the aggregated results.

Intra- and Inter-continent Generalization. The diago-267 nal values represent cases where the model was trained 268 and evaluated on data from the same continent, consistently 269 yielding the highest scores. Notably, Asia achieved the 270 best performance (90.8), followed closely by Europe and 271 the Americas (90.5), indicating strong regional specializa-272 tion. Furthermore, the off-diagonal values reveal models' 273 ability to generalize across regions. Although cross-region 274 scores are generally lower, the model still exhibits reason-275 able transferability. Fine-tuning on Oceania, however, re-276 sulted in the lowest average performance drop, from 85.2 to 277 75.0, suggesting a distinct data distribution for this region. 278

Robustness across Concepts. We grouped the 15 con-279 cepts into 3 main classes: Cultural Heritage and Traditions 280 (CHT), History and Landmarks (HL), and Natural Envi-281 ronment and Local Resources (NELR) using GPT-40. We 282 then conducted training in each group and evaluated on all 283 groups. The specific category mappings are in Appendix B. 284 As shown in Figure 6 (right), models generally perform 285 the best when trained and evaluated within the same cate-286 gory (in-distribution). The model fine-tuned on HL exhibits 287 greater generalization, achieving an average performance of 288 87.8, compared to NELR (84.3). High off-diagonal scores, 289 such as CultureVLM-HL's 87.2 when tested on CHT, indi-290 cate substantial cross-category knowledge transfer, particu-291 larly between culturally related classes. 292

Cross-dataset Generalization. We also evaluated the generalization ability of CultureVLM for cultural reasoning on other datasets. As shown in Figure 8, we tested LLaVA-1.5-7B on CVQA [50] before and after fine-tuning. It is evident that our CultureVLM achieves improvements across most countries and cultures (7% improvement on average), which can be attributed to the comprehensive coverage of our dataset across global nations and cultures, indicating the potential of CultureVerse for cultural research.

Catastrophic Forgetting. Through evaluations on the common VQA benchmarks, we have also found that our method possesses a strong ability to resist catastrophic forgetting. This further demonstrates the strong robustness of our CultureVLM series of models. More details are in E.3.

4. Conclusion

We constructed CultureVerse, a comprehensive multimodal 308 benchmark to characterize the multicultural understanding 309 of VLMs. Extensive evaluation shows significant perfor-310 mance disparities across regions and tasks, highlighting 311 VLMs' strong biases towards Western cultural contexts and 312 their weaker performance in underrepresented regions like 313 Africa and Asia. Using supervised fine-tuning in Culture-314 Verse, we demonstrated effective enhancements in cultural 315 perception and cross-cultural generalization. Our findings 316 underscore the importance of culturally diverse training 317 data and provide actionable insights to improve VLMs. 318

330

331

332

333

350

351

352

353

354

355

356

357

375

376

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

319 References

- [1] Mohammad Amin Abbasi, Arash Ghafouri, Mahdi Firouz mandi, Hassan Naderi, and Behrouz Minaei Bidgoli. Per sianllama: Towards building first persian large language
 model. *arXiv:2312.15713*, 2023. 9
- Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti
 Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach,
 Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3
 technical report: A highly capable language model locally
 on your phone. *arXiv:2404.14219*, 2024. 10
 - [3] Muhammad Farid Adilazuarda, Sagnik Mukherjee, Pradhyumna Lavania, Siddhant Singh, Ashutosh Dwivedi, Alham Fikri Aji, Jacki O'Neill, Ashutosh Modi, and Monojit Choudhury. Towards measuring and modeling" culture" in llms: A survey. arXiv:2403.15412, 2024. 1
- [4] Vibhor Agarwal, Yiqiao Jin, Mohit Chandra, Munmun
 De Choudhury, Srijan Kumar, and Nishanth Sastry. Medhalu: Hallucinations in responses to healthcare queries by
 large language models. *arXiv*:2409.19492, 2024. 13
- [5] Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna,
 Baptiste Bout, Devendra Chaplot, Jessica Chudnovsky,
 Diogo Costa, Baudouin De Monicault, Saurabh Garg,
 Theophile Gervet, et al. Pixtral 12b. *arXiv:2410.07073*,
 2024. 10
- **343** [6] Anthropic. Claude 3.5 sonnet, 2024. 1
- [7] Arjun Appadurai. Modernity at large: Cultural dimensions
 of globalization. *U of Minnesota P*, 1996. 17
- 346 [8] Zechen Bai, Pichao Wang, Tianjun Xiao, Tong He, Zongbo
 347 Han, Zheng Zhang, and Mike Zheng Shou. Halluci348 nation of multimodal large language models: A survey.
 349 arXiv:2404.18930, 2024. 13
 - [9] Mehar Bhatia, Sahithya Ravi, Aditya Chinchure, Eunjeong Hwang, and Vered Shwartz. From local concepts to universals: Evaluating the multicultural understanding of visionlanguage models. arXiv:2407.00263, 2024. 9, 10
 - [10] Y Cao, L Zhou, S Lee, L Cabello, M Chen, and D Hershcovich. Assessing cross-cultural alignment between chatgpt and human societies: An empirical study. arxiv. *Preprint posted online on March*, 31, 2023. 8
- [11] Alex J Chan, José Luis Redondo García, Fabrizio Silvestri,
 Colm O'Donnel, and Konstantina Palla. Harmonizing global
 voices: Culturally-aware models for enhanced content mod eration. *arXiv:2312.02401*, 2023. 9
- [12] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo
 Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou
 Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao,
 and Jifeng Dai. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *arXiv:2312.14238*, 2023. 10
- [13] Yu Ying Chiu, Liwei Jiang, Bill Yuchen Lin, Chan Young
 Park, Shuyue Stella Li, Sahithya Ravi, Mehar Bhatia, Maria
 Antoniak, Yulia Tsvetkov, Vered Shwartz, et al. Culturalbench: a robust, diverse and challenging benchmark
 on measuring the (lack of) cultural knowledge of llms. *arXiv:2410.02677*, 2024. 1, 10
- 374 [14] LMDeploy Contributors. Lmdeploy: A toolkit for com-

pressing, deploying, and serving llm. https://github. com/InternLM/lmdeploy, 2023. 10

- [15] Chengyuan Deng, Yiqun Duan, Xin Jin, Heng Chang, Yijun Tian, Han Liu, Henry Peng Zou, Yiqiao Jin, Yijia Xiao, Yichen Wang, et al. Deconstructing the ethics of large language models from long-standing issues to new-emerging dilemmas. arXiv:2406.05392, 2024. 1
- [16] Yi Fung, Ruining Zhao, Jae Doo, Chenkai Sun, and Heng Ji. Massively multi-cultural knowledge acquisition & Im benchmarking. arXiv:2402.09369, 2024. 2, 8
- [17] Yi R Fung, Tuhin Chakraborty, Hao Guo, Owen Rambow, Smaranda Muresan, and Heng Ji. Normsage: Multi-lingual multi-cultural norm discovery from conversations on-the-fly. *arXiv*:2210.08604, 2022. 8
- [18] Michael Minkov Geert Hofstede, Gert Jan Hofstede. Cultures and Organizations: Software of the Mind, Third Edition. McGraw Hill Professional, https://books.google.co.uk/books?id=7bYWmwEACAAJ, 2010. 8
- [19] Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui, Jie Tang, Jing Zhang, Juanzi Li, Lei Zhao, Lindong Wu, Lucen Zhong, Mingdao Liu, Minlie Huang, Peng Zhang, Qinkai Zheng, Rui Lu, Shuaiqi Duan, Shudan Zhang, Shulin Cao, Shuxun Yang, Weng Lam Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, Xiaohan Zhang, Xiaotao Gu, Xin Lv, Xinghan Liu, Xinyi Liu, Xinyue Yang, Xixuan Song, Xunkai Zhang, Yifan An, Yifan Xu, Yilin Niu, Yuantao Yang, Yueyan Li, Yushi Bai, Yuxiao Dong, Zehan Qi, Zhaoyu Wang, Zhen Yang, Zhengxiao Du, Zhenyu Hou, and Zihan Wang. Chatglm: A family of large language models from glm-130b to glm-4 all tools, 2024. 10
- [20] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-40 system card. *arXiv:2410.21276*, 2024. 1
- [21] Yiqiao Jin, Mohit Chandra, Gaurav Verma, Yibo Hu, Munmun De Choudhury, and Srijan Kumar. Better to ask in english: Cross-lingual evaluation of large language models for healthcare queries. In *Web Conference*, pages 2627–2638, 2024. 1
- [22] Yiqiao Jin, Minje Choi, Gaurav Verma, Jindong Wang, and Srijan Kumar. Mm-soc: Benchmarking multimodal large language models in social media platforms. In ACL, 2024. 10
- [23] Rebecca L Johnson, Giada Pistilli, Natalia Menédez-González, Leslye Denisse Dias Duran, Enrico Panai, Julija Kalpokiene, and Donald Jay Bertulfo. The ghost in the machine has an american accent: value conflict in gpt-3. arXiv:2203.07785, 2022. 8
- [24] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *PNAS*, 114(13):3521–3526, 2017. 14
- [25] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng,
 Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao
 432

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

Zhang, and Ion Stoica. Efficient memory management
for large language model serving with pagedattention. In *SIGOPS*, 2023. 11

- [26] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li,
 Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv:2408.03326*, 2024. 10
- [27] Cheng Li, Mengzhou Chen, Jindong Wang, Sunayana
 Sitaram, and Xing Xie. Culturellm: Incorporating cultural
 differences into large language models. In *NeurIPS*, 2024.
 1, 8, 9, 17
- [28] Cheng Li, Damien Teney, Linyi Yang, Qingsong Wen, Xing
 Xie, and Jindong Wang. Culturepark: Boosting crosscultural understanding in large language models. In *NeurIPS*,
 2024. 8, 9
- [29] Yen-Ting Lin and Yun-Nung Chen. Taiwan Ilm: Bridging the
 linguistic divide with a culturally aligned language model. *arXiv:2311.17487*, 2023. 9
- [30] Chen Cecilia Liu, Fajri Koto, Timothy Baldwin, and Iryna
 Gurevych. Are multilingual llms culturally-diverse reasoners? an investigation into multicultural proverbs and sayings. *arXiv:2309.08591*, 2023. 8
- [31] Fangyu Liu, Emanuele Bugliarello, Edoardo Maria Ponti,
 Siva Reddy, Nigel Collier, and Desmond Elliott. Visually grounded reasoning across languages and cultures. *arXiv:2109.13238*, 2021. 2, 9
- [32] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee.
 Improved baselines with visual instruction tuning. In *CVPR*, pages 26296–26306, 2024. 10
- [33] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee.
 Visual instruction tuning. *NeurIPS*, 36, 2024. 1, 4, 10
- [34] Reem I Masoud, Ziquan Liu, Martin Ferianc, Philip Treleaven, and Miguel Rodrigues. Cultural alignment in large language models: An explanatory analysis based on hofstede's cultural dimensions. *arXiv:2309.12342*, 2023. 8
- 468 [35] Noha Medhat. 8 normal signs and gestures that can be offensive in the middle east. *StepFeed*, 2015. Retrieved 25
 470 February 2019. 1
- 471 [36] Meta. Llama 3.2. https://ai.meta.com/blog/llama-3-2472 connect-2024-vision-edge-mobile-devices, 2024. 10
- 473 [37] Cas Mudde. The 2012 stein rokkan lecture: Three decades
 474 of popu list radical right parties in western europe: so what?
 475 In *The Populist Radical Right*, pages 545–558. Routledge,
 476 2016. 8
- [38] Junho Myung, Nayeon Lee, Yi Zhou, Jiho Jin, Rifki Afina
 Putri, Dimosthenis Antypas, Hsuvas Borkakoty, Eunsu Kim,
 Carla Perez-Almendros, Abinew Ali Ayele, et al. Blend: A
 benchmark for llms on everyday knowledge in diverse cultures and languages. *arXiv:2406.09948*, 2024. 17
- [39] Tarek Naous, Michael J Ryan, Anton Lavrouk, Mohit Chandra, and Wei Xu. Readme++: Benchmarking multilingual language models for multi-domain readability assessment. In *EMNLP*, 2023. 1
- [40] Tarek Naous, Michael J Ryan, and Wei Xu. Having beer after
 prayer? measuring cultural bias in large language models. *arXiv:2305.14456*, 2023. 8

- [41] Shravan Nayak, Kanishk Jain, Rabiul Awal, Siva Reddy, Sjoerd van Steenkiste, Lisa Anne Hendricks, Karolina Stańczak, and Aishwarya Agrawal. Benchmarking vision language models for cultural understanding. *arXiv:2407.10920*, 2024. 492
 1, 9
- [42] Tuan-Phong Nguyen, Simon Razniewski, Aparna Varde, and Gerhard Weikum. Extracting cultural commonsense knowledge at scale. In *Web Conference*, pages 1907–1917, 2023.
 8
- [43] Xuan-Phi Nguyen, Wenxuan Zhang, Xin Li, Mahani Aljunied, Qingyu Tan, Liying Cheng, Guanzheng Chen, Yue Deng, Sen Yang, Chaoqun Liu, et al. Seallms–large language models for southeast asia. arXiv:2312.00738, 2023. 9
- [44] Sejoon Oh, Yiqiao Jin, Megha Sharma, Donghyun Kim, Eric Ma, Gaurav Verma, and Srijan Kumar. Uniguard: Towards universal safety guardrails for jailbreak attacks on multimodal large language models. arXiv preprint arXiv:2411.01703, 2024. 2
- [45] OpenAI. Gpt-40, 2024. 10
- [46] OpenAI. Gpt-4v, 2024. 1
- [47] Kunat Pipatanakul, Phatrasek Jirabovonvisut, Potsawee Manakul, Sittipong Sripaisarnmongkol, Ruangsak Patomwong, Pathomporn Chokchainant, and Kasima Tharnpipitchai. Typhoon: Thai large language models. *arXiv:2312.13951*, 2023. 9
- [48] Ramon Pires, Hugo Abonizio, Thales Sales Almeida, and Rodrigo Nogueira. Sabiá: Portuguese large language models. In *Brazilian Conference on Intelligent Systems*, pages 226–240. Springer, 2023. 9
- [49] Abhinav Rao, Akhila Yerukola, Vishwa Shah, Katharina Reinecke, and Maarten Sap. Normad: A benchmark for measuring the cultural adaptability of large language models. arXiv:2404.12464, 2024. 8
- [50] David Romero, Chenyang Lyu, Haryo Akbarianto Wibowo, Teresa Lynn, Injy Hamed, Aditya Nanda Kishore, Aishik Mandal, Alina Dragonetti, Artem Abzaliev, Atnafu Lambebo Tonja, et al. Cvqa: Culturally-diverse multilingual visual question answering benchmark. *arXiv:2406.05967*, 2024. 1, 2, 4, 9, 10, 16
- [51] Tanik Saikh, Tirthankar Ghosal, Amish Mittal, Asif Ekbal, and Pushpak Bhattacharyya. Scienceqa: a novel resource for question answering on scholarly articles. *International Journal on Digital Libraries*, 23(3):289–301, 2022. 14
- [52] Weiyan Shi, Ryan Li, Yutong Zhang, Caleb Ziems, Raya Horesh, Rogério Abreu de Paula, Diyi Yang, et al. Culturebank: An online community-driven knowledge base towards culturally aware language technologies. *arXiv:2404.15238*, 2024. 1, 8
- [53] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *CVPR*, pages 8317–8326, 2019. 14
- [54] World Values Survey. World values survey. https://www.worldvaluessurvey.org/wvs.jsp, 2022. 8
- [55] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui
 Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al.

546 Gemini: a family of highly capable multimodal models.
547 *arXiv:2312.11805*, 2023. 1, 10

- 548 [56] SM Tonmoy, SM Zaman, Vinija Jain, Anku Rani, Vipula
 549 Rawte, Aman Chadha, and Amitava Das. A comprehensive
 550 survey of hallucination mitigation techniques in large lan551 guage models. arXiv:2401.01313, 2024. 13
- [57] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan,
 Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin
 Ge, et al. Qwen2-vl: Enhancing vision-language model's
 perception of the world at any resolution. *arXiv:2409.12191*,
 2024. 1, 10
- [58] Wenxuan Wang, Wenxiang Jiao, Jingyuan Huang, Ruyi Dai,
 Jen-tse Huang, Zhaopeng Tu, and Michael R Lyu. Not all
 countries celebrate thanksgiving: On the cultural dominance
 in large language models. arXiv:2310.12481, 2023. 8
- [59] Genta Indra Winata, Frederikus Hudi, Patrick Amadeus
 Irawan, David Anugraha, Rifki Afina Putri, Yutong Wang,
 Adam Nohejl, Ubaidillah Ariq Prathama, Nedjma Ousidhoum, Afifa Amriani, et al. Worldcuisines: A massive-scale
 benchmark for multilingual and multicultural visual question
 answering on global cuisines. *arXiv:2410.12705*, 2024. 1
- 567 [60] Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui,
 568 Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui
 569 He, et al. Minicpm-v: A gpt-4v level mllm on your phone.
 570 arXiv:2408.01800, 2024. 10
- [61] Alex Young. Western theory, global world: Western bias in international theory. *Harvard International Review*, pages 29–31, 2014. 1
- [62] Wenlong Zhao, Debanjan Mondal, Niket Tandon, Danica
 Dillion, Kurt Gray, and Yuling Gu. Worldvaluesbench:
 A large-scale benchmark dataset for multi-cultural value
 awareness of language models. *arXiv:2404.16308*, 2024. 8

578	Appendix	
579 580	CultureVLM: Characterizing and Improving Cultural Understanding of Vision-Language Models for over 100 Countries	
581	Contents	
582	A Related Work	8
583 584	B 1 Quality Assurance	9 10
585	B.2. Scalability	10
586	B.3. Analysis of CultureVerse	10
587	C Experiment Setup	10
588	D. Detailed Main Results	11
589	D.1. Task Characteristics: Cultural Awareness Outpaces Image and Detail Recognition	11
590	D.2 Regional Disparities: Better Performance in Western Cultures	11
591	D.3 Category Characteristics: Weak Understanding of History and Landmarks	12
592	D.4. Performance Variability from Model Level	12
593	D.5 Direct Answer v.s. Stepwise Reasoning	13
594	E Detailed Fine-tuning Results	13
595	E.1. Impact of Data Sizes	13
596	E.2. Impact of Decoding Temperatures	14
597	E.3. Catastrophic Forgetting	14
598	F. Details of Human Annotation	15
599	F.1. Statistics of Human Annotators and the Process	15
600	F.2. Accuracy of Human Annotation	16
601	G Limitations	16
602	H Case Study	17
603	I. Prompt List	19

604 A. Related Work

Cultural Bias in LLMs and VLMs. Recent research has increasingly focused on cultural biases present in large language 605 models (LLMs). Johnson et al. [23] investigated conflicts between model outputs and input values and found that GPT-3's 606 responses often aligned more closely with dominant U.S. cultural norms. Similarly, Naous et al. [40] observed a bias toward 607 Western cultural perspectives in models processing Arabic text. The Cultural Alignment Test (CAT), based on Hofstede's 608 609 cultural dimensions framework [18], was used to evaluate the cultural alignment of models like ChatGPT and Bard across various regions, showing that GPT-4 exhibited the strongest alignment with U.S. values [34]. Additionally, Cao et al. [10] 610 found that, while ChatGPT was well-aligned with American cultural values, it struggled to represent other cultures accurately, 611 especially when responding to English prompts. Liu et al. [30] further reported that multilingual LLMs showed limited 612 proficiency in reasoning with proverbs and revealed a "culture gap" in translation tasks. 613

Datasets and Models for Cultural Understanding. Most of the research used existing datasets as cultural datasets. Wang et al. [58] introduced a benchmark based on the World Values Survey (WVS) [54] and the Political Culture and Trust (PCT) dataset [37]. Subsequent works include the Cultural Alignment Test [34], NORMSAGE [17], WorldValueBench [62], and NORMAD [49], each drawing on various existing datasets. Other sources include CultureAtlas [16] and MAPS [30], which collected data from Wikimedia, while Candle [42] and CultureBank [52] gathered data from social media platforms, including TikTok and Reddit. In contrast, there is a growing trend toward automatic data augmentation such as [27, 28]. A strand of

Dataset	Country	Concept	Image	Question	Multi.
MaRVL [31]	5	454	4,914	5,670	X
CVQA [50]	28	-	4,560	9,044	×
CulturalVQA [41]	11	-	2,328	2,328	×
GlobalRG [9]	50/15	-	3,591	-	1
CultureVerse	188	11,085	11,085	31,382	1
w/ Train Set	188	19,682	74,959	196,673	1

Table 1. Comparison of various cultural datasets with features. 'Multi.' indicates whether the dataset provides multi-faceted questions rather than just one single type.

research focuses on training cultural-specific LLMs by assembling large-scale pre-training datasets, followed by fine-tuning to enhance alignment [1, 11, 29, 43, 47, 48]. Instead of relying on massive data collection, Li et al. [27, 28] proposed cost-efficient approaches to fine-tuning cultural-specific models by data augmentation. 622

Unlike LLMs, training data are significantly more difficult to obtain for VLMs. The research in VLMs' cultural bias is still preliminary, with most efforts in *manual* data collection [9, 31, 41, 50]. MaRVL [31] introduced a protocol for building an ImageNet-style hierarchy that represents a wider range of languages and cultures. CVQA [50] proposed a culturally diverse multilingual visual question-answering benchmark designed to encompass a wide variety of languages and cultural contexts, engaging native speakers and cultural experts in the data collection process. CulturalVQA [41] developed a visual questionanswering benchmark focused on evaluating VLMs ' understanding of culturally diverse, geographically specific content. GlobalRG [9] presents two challenging tasks: retrieval across cultural universals and culturally specific visual grounding.

However, these datasets are often limited in size, lack sufficient regional and national representation. More critically, at630the model level, there has been no significant effort to develop culturally aware VLMs, posing a major barrier to AI equity631for underrepresented cultures. By systematically constructing datasets and models through *tangible concepts*, we achieve632both *reduced manual effort* and ensured *scalability*. Table 1 shows the key difference between our benchmark and existing633multimodal cultural datasets.634

B. Details of the CultureVerse Dataset

Table 2 shows the definition of cultural concepts in our dataset. Table 10 shows the number of concepts in different countries636from the evaluation set. In total, we have 11,085 evaluation samples and 19,682 training samples from 188 countries/regions.637

Overall Category	Category	Example Description		
	Festivals	Unique festival celebration scenes		
	Traditional Clothing	Ethnic clothing, festival attire		
Cultural Heritage and	Handicrafts and Artifacts	Ethnic handicrafts, traditional handmade items		
Traditions	Music and Dance	Traditional musical instruments, dance scenes		
	Religion and Belief	Temples and churches, religious ceremonies		
	Entertainment and Performing Arts	Theater performances, street performers		
	Famous Landmarks	Famous historical sites, buildings		
History and Landmarks	Historical Artifacts	Museum collections, ancient relics		
HISTOLY and Landmarks	Historical Figures	Portraits of historical figures, statues		
	Architectural Styles	Traditional architecture, modern landmark buildings		
	Food	Local specialty dishes, traditional festival foods		
Natural Environment	Plants	Unique flowers, crops in a certain area		
and Local Pasources	Animals	Unique wild animals, livestock in a certain area		
and Local Resources	Natural Scenery and Ecosystems	Unique natural landscapes, ecological reserves		
	Markets and Shopping Traditions	Local markets, specialty shops		

Table 2. Collected concepts, their overall categories, and descriptions in the dataset.

635

CVPR

#00006

638 **B.1.** Quality Assurance

To ensure the integrity of the dataset, we conduct a comprehensive manual quality check on each cultural concept, along with its corresponding images and questions. Specifically, we employ human annotators (details are shown in Appendix F) to inspect three main components:

- Image-Concept Alignment. We assessed whether each cultural concept accurately represents the culture of its respective country or region and is either unique to or widely recognized within that. We started with frequency analysis and leverage GPT-40 for preliminary screening, effectively filtering out less desirable data and significantly reducing the manual review workload.
- **Image Quality Check.** We checked the image quality and ensured that the cultural concept is accurately presented in the image.
- Question & Answer Validation. We verified that all three generated questions are reasonable, clear, logically sound, and have a single correct answer. Annotators refined the questions and answer options by removing redundant information and resolving any ambiguities to maintain clarity and accuracy.
- Following the quality assurance process, we utilized human annotations for the evaluation set of CultureVerse while applying the automated annotation pipeline to the larger training set. With over 98% of the evaluation set samples correctly annotated by the automated process, we conclude that the pipeline is highly effective. Any remaining erroneous or challenging samples that could not be refined were filtered out to maintain the dataset's high quality. Additional details on annotator accuracy are in Appendix F.2.

656 B.2. Scalability

Our approach to constructing multimodal cultural datasets is notably more scalable and comprehensive than existing methods. Existing ones mostly rely on manual efforts to search for cultural concepts, retrieve images, and formulate questions, significantly increasing human efforts beyond quality check [13, 22, 50]. This manual process typically results in limited or biased coverage due to the limited scope of cultural concepts explored. For example, some datasets [9] cover only a dozen countries, with approximately 40 cultural concepts for each country.

In contrast, our work is the first to advance the number of specific cultural concepts to the scale of tens of thousands and to extend coverage to approximately two hundred countries/regions, as shown in Table 10. Additionally, our dataset construction process allows for further large-scale expansion, including the retrieval of more images and the synthesis of various types of QAs, such as open-ended QAs, multiple-choice QAs, and reasoning QAs. Our work provides a potential data source for enhancing the multicultural knowledge of VLM models.

667 B.3. Analysis of CultureVerse

668 Figure 2(a) illustrates the distribution of three tasks (Section 2.2), 5 continents (188 countries, with North and South America combined into America) and 15 cultural topics in CultureVerse. Since cultural concepts are collected at the country level, re-669 gions with more countries, such as Asia and Europe, naturally yield larger datasets. Detailed counts of countries and concepts 670 are provided in the Appendix 10. Detailed statistics of the concepts are provided in Table 1. Compared to recent multi-model 671 672 culture datasets, Culture Verse is driven by *tangible*, *presentable cultural concepts*, achieving an order of magnitude increase 673 in the number of regions, images, and questions. This expansion advances multimodal and multicultural research beyond a limited set of countries, moving toward truly global, inclusive cultural integration. Figure 2(b) shows examples of three 674 questions associated with one concept, clearly demonstrating that different question types aim to evaluate different abilities. 675

676 C. Experiment Setup

677 Data Split. To ensure robust evaluation, we partition CultureVerse into training/test sets for all countries/regions, allowing
 678 us to assess transferability between regions. We select more common cultural concepts from the entire dataset for the test set,
 679 which underwent manual quality checks, while the training set includes all cultural concepts. We ensured that the images in
 680 training and test sets did not overlap to prevent data leakage.

Evaluation Models. We conduct evaluations on the following models: (a) open-source models including LLaVA-1.5 [32],
LLaVA-1.6-Mistral-7B-Instruct [33], LLaVA-OneVision [26], LLaMA-3.2-Vision [36], Qwen2-VL [57], InternVL-2 [12],
Phi-3-Vision [2], MiniCPM-Llama3-V-2.5 [60], GLM-4V [19], Pixtral 12B [5]; (b) proprietary models such as GPT-40 [45],

- 684 Gemini-1.5-Pro [55].
- **Evaluation Setup.** We report accuracy, in line with previous works [33]. For all multiple-choice questions, we employ greedy search decoding for deterministic predictions. For LLaVA-1.5, LLaMA-3.2-Vision and InternVL-2, we use Imdeploy [14]

707

708

720

for inference acceleration. For other models, we use vllm [25] for the acceleration of inference. For models before and 687 after fine-tuning, we use the same acceleration toolkit to prevent potential impact. The number of questions differs across 688 the three tasks. This is because generating questions for cultural knowledge and scenario reasoning is more complex, and 689 in some cases, GPT-40 refused to provide answers, making it impossible to generate valid questions. For image recognition 690 questions, we directly use the questions and options as prompts. For cultural knowledge and scene understanding questions, 691 we employ stepwise reasoning prompts to facilitate the reasoning explanation. The prompts are available in the Appendix I. 692 For all proprietary models, we utilize the default hyper-parameters, setting the temperature to 0 and the max tokens to 1,024. 693 For all open-source models, *do_sample* is set to False, *max_gen_len* is set to 512, and the temperature is set to 0.01. 694 Training Setup. We use the official train ing scripts of LLaVA², Phi³, and LLaMA⁴ for model training, largely adhering to 695 the original hyperparameters, except for appropriately adjusting the batch size to accommodate the GPU memory capacity. 696 For LLaVA-1.5, a learning rate of 2×10^{-5} is used, with no weight decay applied (0.0). The learning rate followed a cosine 697 schedule, gradually increasing during the initial phase with a warmup ratio of 0.03. For Phi-3, we use a learning rate of 698 4×10^{-5} and a weight decay of 0.01. A linear learning rate scheduler is utilized, with 50 warmup steps to stabilize the early 699 training stage. For LLaMA-3.2, fine-tuning is conducted using a learning rate of 1×10^{-5} with no weight decay (0.0). A 700 multiplicative learning rate decay is applied after each epoch, with a gamma value of 0.85. The batch sizes are set to 64, 16 701 and 32 respectively. All models are trained for one epoch on the training set and fully fine-tuned on 4×A100 80GB GPUs. 702 For the training data, although we do not conduct large-scale human annotation, we synthesize the data using only concepts 703 that passed either GPT-40 or human quality assurance, significantly improving the accuracy of the dataset. The prompts used 704 for GPT quality check can be found in the Appendix I. 705

D. Detailed Main Results

Detailed results on different tasks, continents, and cultural categories can be found in Table 3.

D.1. Task Characteristics: Cultural Awareness Outpaces Image and Detail Recognition

From the task perspective, we observe that image recognition and cultural knowledge questions pose challenges comparable 709 for VLMs. *Image recognition* tests VLMs' ability to identify culturally specific objects or concepts, which relies heavily 710 on diverse and relevant image data. For instance, recognizing traditional foods like kimchi from Korea, or regional attire 711 such as a sari from India, requires the model to have encountered similar image-text pairs in its training data. In contrast, 712 cultural knowledge questions assess the model's understanding of broader cultural elements based on text-based training. For 713 example, asking about the significance of a festival like *Diwali* or the symbolism of a *red envelope* during Lunar New Year 714 taps into the model's text-based memory, which tends to be richer due to the abundance of internet text data. Interestingly, 715 in the scene understanding task, which integrates images with contextual background (e.g., a Japanese tea ceremony scene 716 or a Brazilian Carnival parade), VLMs tend to generate culturally appropriate responses, avoiding inappropriate or culturally 717 insensitive outputs. This can be attributed to the model's inherent multicultural awareness and its alignment with ethical and 718 harm-reduction training. 719

D.2. Regional Disparities: Better Performance in Western Cultures

A dominant *regional disparity* is observed among all models: VLMs demonstrate the strongest cultural understanding of the 721 Americas, followed by Europe and Oceania. This trend reflects the dominance of English data centered around the Global 722 North, leading to a disproportionate focus on Western cultural content. North America's relatively homogenous cultural 723 landscape, combined with fewer countries, contributes to better model performance. In contrast, Asia and Africa show 724 significantly weaker results, likely due to the scarcity of digitized, English-language data and the high cultural diversity in 725 these regions. For instance, Asia consists of many countries with distinct and complex cultural contexts, such as those from 726 East Asia, South Asia, and Southeast Asia, both within and across nations. Although Asia has the most data in CultureVerse 727 (see Figure 2), the models struggle to capture the intra-regional and inter-regional cultural nuances, resulting in suboptimal 728 performance. 729

²https://github.com/haotian-liu/LLaVA.

³https://github.com/microsoft/Phi-3CookBook.

⁴https://github.com/meta-llama/llama-recipes/blob/main/recipes/quickstart/finetuning/finetune_vision_ model.md.

CVPR 2025 Submission #00006. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.







Figure 4. Results and analysis of CultureVLM by fine-tuning on our CultureVerse.

730 D.3. Category Characteristics: Weak Understanding of History and Landmarks

VLMs generally exhibit weaker recognition and understanding of cultural concepts related to history and landmarks. The
 primary reason is the relatively limited internet data available on historical figures and landmarks. Additionally, recognizing
 a landmark typically requires training data that includes images from multiple perspectives to form a comprehensive, three dimensional understanding.

735 D.4. Performance Variability from Model Level

Proprietary models continue to outperform open-source counterparts, with GPT-40 achieving the best results. Although larger
 models tend to demonstrate better performance, size alone is not the determining factor. For example, the size variations of
 LLaVA-1.5 and Qwen2-VL show similar performance. Cultural knowledge often resides in model's memory, an aspect
 overlooked in VLMs development. Thus, smaller models (e.g. Phi-3-Vision) with comparable training data can already

741

exhibit strong cultural understanding when using similar training data as larger models.

D.5. Direct Answer v.s. Stepwise Reasoning

For image recognition tasks, we compare two prompt methods: 1) the model directly identifies and outputs the cultural concept, and 2) first provides a detailed description of the image and then analyzes and compares the options to reach a final answer. As shown in Figure 4a, we find that stepwise reasoning does not improve cultural recognition and, in most cases, significantly *impair* performance. Upon analyzing the answers, we observe that VLMs frequently exhibit hallucinations [4, 56] during step-by-step explanations, as shown in the example in Figure 9. This poor robustness suggests that while VLMs rate understanding of the details and components that make up those concepts. 748

		Task				Contine	ent			Categor	у
Model	Image Recognition	Cultural Knowledge	Scene Reasoning	Africa	Asia	Europe	America	Oceania	СНТ	HL	NELR
			Open-Sc	ource Mod	lels						
LLAVA-v1.57B	51.19	48.13	72.44	55.62	53.78	57.24	63.27	60.41	58.73	56.27	58.75
LLAVA-v1.5 13B	53.05	38.51	61.40	49.83	47.23	51.52	57.53	54.51	52.21	50.28	53.24
LLAVA-NEXT-v1.67B	44.48	57.03	80.61	58.80	57.67	60.03	65.61	63.72	62.78	58.77	63.22
LLAVA-ONEVISION 0.5B	45.38	36.40	65.33	47.80	48.24	48.23	51.48	53.59	51.19	48.38	48.79
LLAVA-ONEVISION 7B	64.50	61.98	78.94	66.00	66.53	68.59	72.20	70.63	68.50	67.98	70.17
InternVL2 2B	41.86	40.99	64.22	48.10	47.01	48.73	52.67	50.46	50.22	48.15	50.24
InternVL2 8B	58.40	65.34	85.43	65.58	66.89	70.32	73.94	71.92	69.47	69.34	70.05
Phi-3-vision 4B	54.85	56.43	79.09	60.60	60.99	63.21	67.78	68.69	64.22	62.40	65.73
MiniCPM-V-2_6 8B	65.07	71.69	88.00	71.52	72.59	76.10	77.60	74.86	74.77	74.62	74.93
GLM-4V 9B	68.59	66.03	82.25	69.38	69.34	73.41	76.67	73.20	71.86	72.45	72.02
Pixtral 12B	64.43	72.56	89.82	71.79	73.04	76.12	79.41	77.44	75.56	74.39	79.02
Qwen2-VL 7B	75.51	71.90	77.96	73.59	72.73	75.71	78.93	77.90	75.05	74.57	77.78
Qwen2-VL 72B	72.73	73.20	83.05	76.62	73.87	75.71	80.66	80.11	78.30	74.69	79.84
LLaMA-3.2-Vision 11B	71.87	73.72	87.86	74.55	76.62	77.57	80.90	79.28	78.12	77.24	79.04
			Proprie	tary Mod	els						
Gemini-1.5-Pro	81.88	87.03	93.60	84.91	86.14	88.06	89.62	87.08	86.50	87.43	88.36
GPT-40	84.67	90.86	96.16	88.69	90.22	90.91	90.92	89.41	91.53	89.50	92.58

Table 3. Zero-shot accuracy on open-source and proprietary models across three tasks, five continents, and three categories. CHT denotes *Cultural Heritage and Traditions*; HL denotes *History and Landmarks*; and NELR denotes *Natural Environment and Local Resources*.

E. Detailed Fine-tuning Results

The detailed results of the fine-tuned models are shown in Table 4. Detailed results on the generalization of the fine-tuned model in different regions and for different categories can be found in Table 5 and Table 6. Detailed results of the models before and after fine-tuning on the general VQA benchmark are shown in Table 7.

		Task				Contine	ent			Categor	у
Model	Image Recognition	Cultural Knowledge	Scene Reasoning	Africa	Asia	Europe	America	Oceania	СНТ	HL	NELR
LLAVA-v1.57B	88.03	90.87	95.66	91.68	91.61	91.02	91.91	90.61	92.58	90.98	91.70
LLAVA-v1.5 13B	89.72	92.20	96.09	92.68	92.37	92.59	93.05	92.73	93.17	92.40	92.60
Phi-3-vision 4B	87.53	90.84	95.77	90.80	91.56	91.41	90.91	91.16	92.83	90.50	92.35
LLaMA-3.2-Vision 11B	89.13	91.49	96.20	91.99	92.24	91.82	93.08	91.34	93.20	91.78	92.50

Table 4. Performance of fine-tuned models across three tasks, five continents, and three categories. CHT denotes *Cultural Heritage and Traditions*; HL denotes *History and Landmarks*; and NELR denotes *Natural Environment and Local Resources*.

E.1. Impact of Data Sizes

We vary the number of fine-tuning examples within [5%, 20%, 50%, 75%, 100%] to examine the effects of training data sizes 754 on the final results. As shown in Figure 4c and 4d, the model's performance decreases as the training data is reduced. 755

13

749

752

753

CVPR 2025 Submission #00006. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.



Figure 5. Performance of open-source and proprietary models on general VQA datasets. *Base* refers to the original model, while *Finetuned* refers to the model adapted using our CultureVerse. The comparable performance across both versions suggests that finetuning on our dataset preserves the models' natural language understanding and commonsense reasoning abilities.



Figure 6. Generalization and Robustness. Left: Performance of CultureVLM (y-axis) evaluated across data from different continents (x-axis). CultureVLM achieves the highest performance for in-distribution settings, while still demonstrating strong generalizability for out-of-domain settings. **Right**: CultureVLM fine-tuned with data under different categories of CultureVerse (x-axis) and evaluated across various categories (y-axis). CHT denotes *Cultural Heritage and Traditions*; HL denotes *History and Landmarks*; and NELR denotes *Natural Environment and Local Resources*.

However, the decline is minimal, indicating that even a small amount of training data can effectively enhance the model'smulticultural awareness.

758 E.2. Impact of Decoding Temperatures

We evaluate the performance of the original models and CultureVLMs under different temperature and decoding settings, as shown in Figure 7. It can be observed that VLMs perform better when the temperature is lower and decoding diversity is reduced. However, when the temperature reaches 1.0, there is a noticeable and expected decline in performance.

762 E.3. Catastrophic Forgetting

763 Catastrophic forgetting [24] refers to the phenomenon where a model loses previously learned knowledge when trained 764 on new information, especially when the new data diverges significantly from the pretraining data. This can be especially 765 problematic in cultural knowledge acquisition, as it may cause the model to compromise essential commonsense knowledge 766 in favor of culture-specific details.

To assess this, we evaluate the models on standard VQA benchmarks, including ScienceQA [51] and TextVQA [53] to determine if the process of acquiring cultural knowledge inadvertently diminishes their grasp of general commonsense concepts. The results in Figure 5 reveal that our fine-tuned CultureVLM merely influences general VLM tasks, indicating the versatility of the solution.

772



Figure 7. Impact of different decoding temperatures on performance

Task	Model	Total	Africa	Asia	Europe	America	Oceania
	LLaVA-v1.5-7B-Africa	72.73	85.79	72.69	70.54	71.37	69.74
	LLaVA-v1.5-7B-America	75.84	75.35	73.24	72.52	86.00	76.05
Image Recognition	LLaVA-v1.5-7B-Asia	80.20	75.78	88.62	75.04	75.84	75.26
	LLaVA-v1.5-7B-Europe	79.34	77.61	75.64	86.89	75.62	76.05
	LLaVA-v1.5-7B-Oceania	64.51	64.05	63.14	63.60	65.28	83.95
	LLaVA-v1.5-7B-Africa	78.54	86.77	75.82	78.28	80.38	79.42
	LLaVA-v1.5-7B-America	81.85	79.29	78.33	80.45	90.78	85.22
Cultural Knowledge	LLaVA-v1.5-7B-Asia	84.42	81.25	88.76	81.27	82.90	83.77
	LLaVA-v1.5-7B-Europe	84.12	80.52	80.20	90.43	83.15	81.74
	LLaVA-v1.5-7B-Oceania	73.78	73.16	72.38	73.91	75.03	81.16
	LLaVA-v1.5-7B-Africa	91.29	93.63	91.06	90.87	91.41	91.41
	LLaVA-v1.5-7B-America	92.70	93.40	92.24	91.85	94.81	91.69
Scene Understanding	LLaVA-v1.5-7B-Asia	94.16	94.44	95.11	92.92	94.52	92.52
	LLaVA-v1.5-7B-Europe	93.26	92.71	92.64	94.08	93.45	92.52
	LLaVA-v1.5-7B-Oceania	87.54	86.57	88.08	86.55	88.11	90.58

Task	Model	Total	Cultural Heritage and Traditions	History and Landmarks	Natural Environment and Local Resources
	LLaVA-v1.5-7B-CHT	77.91	88.58	74.36	77.48
Image Recognition	LLaVA-v1.5-7B-HL	84.76	81.57	86.89	80.25
	LLaVA-v1.5-7B-NELR	73.46	74.89	69.18	89.88
	LLaVA-v1.5-7B-CHT	83.29	88.53	81.71	82.16
Cultural Knowledge	LLaVA-v1.5-7B-HL	87.99	84.91	90.58	82.03
	LLaVA-v1.5-7B-NELR	81.17	80.71	79.91	86.84
	LLaVA-v1.5-7B-CHT	92.96	95.27	92.45	91.73
Scene Understanding	LLaVA-v1.5-7B-HL	94.77	95.23	94.90	93.62
	LLaVA-v1.5-7B-NELR	91.85	92.58	91.16	93.62

Table 5. Accuracy across different continents for each fine-tuned model.

Table 6. Accuracy across different categories for each fine-tuned model.

F. Details of Human Annotation

F.1. Statistics of Human Annotators and the Process

Table 8 shows the statistics of human annotators in our study. In total, through the contractor company, we hired 10 expert773annotators whose ages are between 20 and 36 with at least a bachelor's degree. Most of them are within the non-AI areas774such as education, specific languages, and history. When assigning the annotation job, we asked each annotator to label the775correctness, consistency, and relatedness of our questions and answers. Specifically, correctness refers to the correctness of776the generated questions and answers, consistency refers to the consistency between the questions, and the concepts,777and relatedness aims to make sure that the concepts and generated questions are related to each other. We asked the annotators778not only to make judgement based on their experience but also to manually check the results via Google search and other779

CVPR 2025 Submission #00006. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

Dataset	Model	Base	Finetuned
ScienceQA	LLaVA 7B	70.12	69.28
	LLaVA 13B	74.91	75.04
	LLAMA 3.2 11B	88.97	89.30
TextVQA	LLaVA 7B FT	58.32	57.89
	LLaVA 13B FT	61.18	61.13
	LLAMA 3.2 11B	71.34	70.67

Table 7. Performance of models before and after fine-tuning on general VQA datasets. *Base* refers to the original model, while *Finetuned* refers to the model adapted using our CultureVerse. The comparable performance across both versions suggests that finetuning on our dataset preserves the models' natural language understanding and commonsense reasoning abilities.

Age	%	Degree	%	Major
20-25	50%	Bachelor	50%	Education; Specific languages; Computer science;
26-36	50%	Master	50%	Communication; Public relation; History

Table 8.	Statistics of	the human	annotators to	validate	CultureVerse.
----------	---------------	-----------	---------------	----------	---------------

search engines. Each instance is labeled by two experts and then verified by another to ensure correctness. All annotation
 operations are performed following local laws and regulations to ensure fairness, equity, and accountability.

F.2. Accuracy of Human Annotation

Table 9 shows the precision of human annotators in our generated questions. We then filter out all the wrong questions and
 only retain the correct ones. It is surprising that automatically generated questions can achieve high accuracy, indicating the
 promising future of the adoption of advanced AI models like GPT-40 for data collection and annotation.

Check Item	Accuracy
Concept-Region Alignment	98.25%
Concept-Image Alignment	99.49%
Image Recognition Question-Answer Correctness	98.61%
Cultural Knowledge Question-Answer Correctness	96.52%
Scene Understanding Question-Answer Correctness	93.18%

Table 9. Accuracy of CultureVerse based on the human annotations.



Figure 8. Results on CVQA [50].

Figure 9. Case study on the effect of fine-tuning and prompt variations.

786 G. Limitations

First, while we use languages as proxies for cultural boundaries in the development of CultureVLM, we acknowledge that
 language alone does not capture the full complexity of culture. This simplification was made to address challenges in defining

cultural contexts, inspired by previous research [7, 27, 38]. Our pipeline and approach remain flexible to incorporate addi-789 tional languages and cultures. Second, we acknowledge that the current dataset lacks multilingual support. Most foundational 790 models currently exhibit weak multilingual capabilities, so fine-tuning on multilingual cultural data is less effective than on 791 English data [27]. Therefore, the English data provided currently is also timely and valuable to our community. Third, due 792 to resource constraint cost, we limit our fine-tuning experiments to the current set of models. Evaluating a wider range of 793 models could yield further insight. The current CultureVerse only contains multiple-choice questions. Exploring open-ended 794 questions could also offer additional avenues for assessment. Our data and models will be released to the public after ethical 795 reviews. 796

H. Case Study

797

802

Figure 9 shows the responses of LLaVA and CultureVLM. We incorporate extensive explanations into the training data, enriching CultureVLM with substantial knowledge that enhances its cultural recognition and understanding ability. 799

Below, we present examples from CultureVerse representing three different countries. Each example includes a cultural concept, country, category, image recognition question, cultural knowledge question, and scene understanding question.



Concept: Florence Cathedral Country: Italy Category: Famous Landmarks

Image Recognition:

Question: What famous cathedral is shown in this image? Options: (A) St. Peter's Basilica (B) Milan Cathedral (C) Florence Cathedral (D) Siena Cathedral Ground truth: (C) Florence Cathedral

Cultural Knowledge:

Question: Which renowned architect was responsible for engineering the innovative dome of the structure shown in the image?
Options: (A) Leon Battista Alberti (B) Filippo Brunelleschi (C) Giorgio Vasari (D) Michelangelo Buonarroti
Ground truth: (B) Filippo Brunelleschi

Scene Understanding:

Question: Imagine you are visiting the city depicted in the image during the Renaissance period. You hear local legends about a stone head on a famous structure there. According to the legend, what was the purpose of this stone head? **Options:** (A) To ward off evil spirits (B) To grant wishes to passersby (C) To commemorate a donor who contributed a bell (D) To mark the entrance for pilgrims **Ground truth:** (C) To commemorate a donor who contributed a bell

Concept: Pha That Luang **Country:** Lao People's Democratic Republic **Category:** Famous Landmarks

Image Recognition:

Question: What famous structure is shown in this image? **Options:** (A) Wat Arun (B) Angkor Wat (C) Pha That Luang (D) Shwedagon Pagoda **Ground truth:** (C) Pha That Luang



Cultural Knowledge:

Question: The structure shown in the image, originally built as a Hindu temple, underwent major reconstruction during the reign of which historical figure when Buddhism gained prominence in the region?

Options: (A) King Jayavarman VII (B) King Setthathirath (C) King Anawrahta (D) King Ramkhamhaeng

Ground truth: (B) King Setthathirath

Scene Understanding:

Question: Imagine you are attending a festival at the site shown in the image, which is considered one of the most significant Buddhist celebrations in the country. During this event, which of the following cultural practices would you most likely observe that emphasizes the religious and national significance of this site?

Options: (A) A parade featuring traditional music and dance (B) A ceremony honoring ancient Hindu deities (C) A cooking contest of traditional Lao dishes (D) A fireworks display celebrating the lunar new year

Ground truth: (A) A parade featuring traditional music and dance

I. Prompt List

806 Concept: {concept} 807 Country: {country} 808 Class: {class_} 809 810 Please determine whether "{concept}" is a kind of {class_} that can reflect {country} 811 culture and whether it can serve as a symbol of {country} culture (unique and very famous 812 813 within {country}, and not commonly seen in other parts of the world). Additionally, the symbol should not be a broad category that includes various specific items 814 , but rather a distinct and indivisible entity, such as a specific dance form, a famous 815 individual's photograph, a renowned landmark, etc. 816 817 818 Here are some counterexamples: - Broad concepts like "fast food" (which includes burgers, fries, etc.), "traditional 819 Chinese instruments" (which include guzheng, erhu, etc.) "Dance" (which include jazz dance, 820 Square dancing, etc.) are not acceptable due to their lack of a unified visual marker. 821 - Concept itself is the wrong word, such as "...", "N/A". 822 - Concept exists in many regions, such as "pork", "duck" or "grapes" (which might be a 823 specialty or staple in certain countries but is quite common in many places. However, 824 specific concepts like "peking duck" or "schweinshaxe" would be correct). 825 826 827 Please provide your short explanation and include your answer (Yes/No) into <<<>>>. For example, if you think "{concept}" is a specific and indivisible symbol of {country} culture, 828 please write <<<Yes>>>. 838 1. Prompt of concept judgement. 831 Please extract cultural elements from the given wikipedia document that can represent { 832 833 country} culture or are very famous in {country}, including the following categories: 834 ### Categories 835 - Food: e.g., local specialty dishes, traditional festival foods. 836 - Plants: e.g., unique flowers, crops in {country}. 837 - Animals: e.g., unique wild animals, livestock in {country}. 838 839 - Famous Landmarks: e.g., famous historical sites, buildings. - Festivals: e.g., unique festival celebration scenes. 840 - Historical Artifacts: e.g., museum collections, ancient relics. 841 - Historical Figures: e.g., portraits of historical figures, statues. 842 - Traditional Clothing: e.g., ethnic clothing, festival attire. 843 - Architectural Styles: e.g., traditional architecture, modern landmark buildings. 844 - Handicrafts and Artifacts: e.g., ethnic handicrafts, traditional handmade items. 845 846 - Music and Dance: e.g., traditional musical instruments, dance scenes. - Religion and Belief: e.g., temples and churches, religious ceremonies. 847 - Natural Scenery and Ecosystems: e.g., unique natural landscapes, ecological reserves. 848 849 - Markets and Shopping Traditions: e.g., local markets, specialty shops. - Entertainment and Performing Arts: e.g., theater performances, street performers. 850 851 Please note that not all categories may be included in the document. Only list the most 852 famous cultural elements, with a total not exceeding 10. For categories without famous 853 elements, use 'NA' to indicate. Directly output in the format "Category: [Element1, Element2 854 , ...]", for example: 855 856 ### Cultural Elements 857 - Food: [Food 1], [Food 2], ... 858 859 - Plants: NA - Music and Dance: [Musical Instrument 1], [Dance Scene 2], ... 860 861

805

	2. Prompt for extracting cultural concept entities from Wikipedia documents.
This is asks th	the {concept} of {country}. Your task is to generate a multiple-choice question the user to identify what is shown in the image. Use the following format for your
- Questi Option 4	.on: [Your Question] Options: (A) [Option 1] (B) [Option 2] (C) [Option 3] (D) [!]
For exam like th	uple, if the image shows a Peking Duck of China, the question and options should in is:
– Questi (B) Peki	on: What traditional dish is shown in this image? Options: (A) Cantonese Roast Dung Duck (C) Sichuan Spicy Duck (D) Nanjing Salted Duck
If the i - Questi Sanxian	mage shows the Erhu of China: .on: What musical instrument is shown in this image? Options: (A) Pipa (B) Erhu ((D) Yangqin
If the i - Questi B) The W	mage shows the White House of US: .on: What famous American building is shown in this image? Options: (A) The Capito White House (C) The Lincoln Memorial (D) The Supreme Court Building
Ensure t are plau enough t potentia	that "{concept}" should be included in one of the options. Ensure that the options isible but only one is the correct answer. The incorrect options should be similat to the correct one to create a challenge, but not so similar that they cause any al ambiguity.
	3. Prompt for generating scene recognition questions.
Please p as: Loc Time per developm }, moder associat	provide a detailed introduction of {concept} of {country}, including information station and Features (where it is found or originates from, what makes it unique); ciod (when it was created or became significant); History (historical background attent, or any significant events); Cultural significance (cultural Context in {countrol and significance); Stories or Legends (Any stories, legends, or folklore attent {concept}).
Use the Introduc	following format for your introduction: ction of {concept}: [Detailed Introduction of {concept}]
Here are	e two examples:
Introduc during staple i	tion of Peking Duck: Peking Duck is a famous Chinese dish that originated in Bei the Imperial era. The dish dates back to the Yuan Dynasty (1271-1368) and became In the Ming Dynasty (1368-1644). Traditionally, Peking Duck is known for its thin
crispy s involves	kin and is served with pancakes, hoisin sauce, and scallions. The preparation inflating the duck to separate the skin from the fat, marinating it, and roastin closed or hung oven. It is considered a national dish of China and a symbol of culinary art.
Chinese	

<pre>4. Prompt for generating the introduction of cultural concept. This public image shows the "[concept]" of (country]. Generate a multiple-choice question based on this image and the introduction of (concept). Provide the correct answer immediately following the question. Baser the question delves into deeper cultural knowledge but does not directly name the { concept). The options should be somewhat confusing to increase the difficulty, but there must be only one correct answer. Users can only answer based on the image, so don't mention my "introduction" or "(concept)" in the question. Use the following format for your generated question: - Question: [Your Question] Options: (A) [Option 1] (B) [Option 2] (C) [Option 3] (D) [Option 4] - Answer: (X) [Option X] Here are two examples: Image: Peking Duck Introduction of Feking Duck: Peking Duck is a famous Chinese dish that originated in Beijing during the Imperial era. The dish dates back to the Yuan Dynasty (1271-1368) and became a staple in the Ming Dynasty (1366-1644). Traditionally, Peking Duck is known for its thin, rrispy skin and is served with pancakes, hoisin sauce, and scallions. The preparation involves inflating the duck to separate the skin from the fat, marinating it, and roasting ti na closed or hung oven. It is considered an ational dish of China and a symbol of Chinese culinary art. - Question: During which dynasty did the dish shown in the image become a staple in the subsine of its country? Options: (A) Tang Dynasty (B) Song Dynasty (C) Ming Dynasty (D) Qing Dynasty - Answer: (C) Ming Dynast</pre>	Now please provide the introduction for {concept}. Introduction of {concept}:					
This public image shows the "{concept}" of {country}. Cenerate a multiple-choice question based on this image and the introduction of {concept}. Provide the correct answer immediately following the question. Snsure the question delves into deeper cultural knowledge but does not directly name the { concept}. The options should be somewhat confusing to increase the difficulty, but there must be only one correct answer. Users can only answer based on the image, so don't mention my "introduction" or "(concept)" in the question. Use the following format for your generated question: - Question: [Your Question] Options: (A) [Option 1] (B) [Option 2] (C) [Option 3] (D) [Dption 4] - Answer: (X) [Option X] dere are two examples: Image: Peking Duck Introduction of Peking Duck: Peking Duck is a famous Chinese dish that originated in Beijing during the Imperial era. The dish dates back to the Yuan Dynasty (1271-1368) and became a staple in the Ming Dynasty (1368-1644). Traditionally, Peking Duck is known for its thin, prispy skin and is served with pancakes, hoisin sauce, and scallions. The preparation involves inflating the duck to separate the skin from the fat, marinating it, and roasting it in a closed or hung oven. It is considered a national dish of China and a symbol of Chinese culinary art. - Question: During which dynasty did the dish shown in the image become a staple in the zuisine of its country? Options: (A) Tang Dynasty (B) Song Dynasty (C) Ming Dynasty (D) Qing Dynasty - Answer: (C) Ming Dynasty Introduction of The White House: The White House, located at 1600 Pennsylvania Avenue NW in Mashington, D.C., Is the official residence and workplace of the President of the United States. Construction began in 1792 and was completed in 1800. The building was designed by Frish-born architect James Hoban In the neoclassical style. It has been the recoldence of svery U.S. president since John Adama. The White House has undergone several renovations and expansions, including the addition of the West Wing and the Oval Of	4. Prompt for generating the introduction of cultural concept.					
<pre>Znsure the question delves into deeper cultural knowledge but does not directly name the (concept). The options should be somewhat confusing to increase the difficulty, but there must be only one correct answer. Users can only answer based on the image, so don't mention any "introduction" or "(concept)" in the question. Use the following format for your generated question: - Question: [Your Question] Options: (A) [Option 1] (B) [Option 2] (C) [Option 3] (D) [Option 4] - Answer: (X) [Option X] Here are two examples: Image: Peking Duck Introduction of Peking Duck: Peking Duck is a famous Chinese dish that originated in Beljing during the Imperial era. The dish dates back to the Yuan Dynasty (1227-1368) and became a staple in the Ming Dynasty (1366-1644). Traditionally, Peking Duck is known for its thin, rrispy skin and is served with pancakes, hoisin sauce, and scallions. The preparation involves inflating the duck to separate the skin from the fat, marinating it, and roasting it in a closed or hung oven. It is considered a national dish of China and a symbol of Chinese culinary art. - Question: During which dynasty did the dish shown in the image become a staple in the puisine of its country? Options: (A) Tang Dynasty (B) Song Dynasty (C) Ming Dynasty (D) Qing Dynasty - Answer: (C) Ming Dynasty Image: The White House Introduction of The White House: The White House, located at 1600 Pennsylvania Avenue NW in Kashington, D.C., is the official residence and workplace of the President of the Dinited States. Construction began in 1792 and was completed in 1800. The building was designed by Irish-born architect James Hoban in the neoclassical style. It has been the residence of avery U.S. president shoe Andama. The White House has undergone several renovations and expansions, including the addition of the West Wing and the Oval Office. It is a symbol of the U.S. government and a site of significant historical events. - Question: Who was the architect responsible for designing the building shown in the image? Op</pre>	This public image shows the "{concept}" of {country}. Generate a multiple-choice question based on this image and the introduction of {concept}. Provide the correct answer immediately following the question.					
<pre>Question: [Your Question] Options: (A) [Option 1] (B) [Option 2] (C) [Option 3] (D) [pption 4] Answer: (X) [Option X] Here are two examples: Image: Peking Duck Introduction of Peking Duck: Peking Duck is a famous Chinese dish that originated in Beijing during the Imperial era. The dish dates back to the Yuan Dynasty (1271-1368) and became a staple in the Ming Dynasty (1368-1644). Traditionally, Peking Duck is known for its thin, rrispy skin and is served with pancakes, hoisin sauce, and scallions. The preparation involves inflating the duck to separate the skin from the fat, marinating it, and roasting it in a closed or hung oven. It is considered a national dish of China and a symbol of Chinese culinary art Question: During which dynasty did the dish shown in the image become a staple in the usisine of its country? Options: (A) Tang Dynasty (B) Song Dynasty (C) Ming Dynasty (D) Qing Dynasty - Answer: (C) Ming Dynasty Image: The White House Introduction of The White House: The White House, located at 1600 Pennsylvania Avenue NW in Asshington, D.C., is the official residence and workplace of the President of the United States. Construction began in 1792 and was completed in 1800. The building was designed by Irish-born architect James Hoban in the moclassical style. It has been the residence of avery U.S. president since John Adams. The White House has undergone several renovations and expansions, including the addition of the West Wing and the Oval Office. It is a symbol of the U.S. government and a site of significant historical events Question: Who was the architect responsible for designing the building shown in the image? Options: (A) James Hoban Now please generate the question for the Image: (concept) of (country) [introduction]</pre>	Ensure the question delves into deeper cultural knowledge but does not directly name the { concept}. The options should be somewhat confusing to increase the difficulty, but there must be only one correct answer. Users can only answer based on the image, so don't mention any "introduction" or "{concept}" in the question. Use the following format for your generated question:					
<pre>Here are two examples: Image: Peking Duck Introduction of Peking Duck: Peking Duck is a famous Chinese dish that originated in Beijing during the Imperial era. The dish dates back to the Yuan Dynasty (1271-1368) and became a staple in the Ming Dynasty (1368-1644). Traditionally, Peking Duck is known for its thin, prispy skin and is served with pancakes, hoisin sauce, and scallions. The preparation involves inflating the duck to separate the skin from the fat, marinating it, and roasting it in a closed or hung oven. It is considered a national dish of China and a symbol of Chinese culinary art. - Question: During which dynasty did the dish shown in the image become a staple in the suisine of its country? Options: (A) Tang Dynasty (B) Song Dynasty (C) Ming Dynasty (D) Qing Dynasty - Answer: (C) Ming Dynasty Image: The White House Introduction of The White House: The White House, located at 1600 Pennsylvania Avenue NW in Washington, D.C., is the official residence and workplace of the President of the United States. Construction began in 1792 and was completed in 1800. The building was designed by Trish-born architect James Hoban in the neoclassical style. It has been the residence of every U.S. president since John Adams. The White House has undergone several renovations and expansions, including the addition of the West Wing and the Oval Office. It is a symbol of the U.S. government and a site of significant historical events. - Question: Who was the architect responsible for designing the building shown in the image? Options: (A) James Hoban Now please generate the question for the Image: {concept} of {country} (introduction)</pre>	- Question: [Your Question] Options: (A) [Option 1] (B) [Option 2] (C) [Option 3] (D) [Option 4] - Answer: (X) [Option X]					
<pre>Image: Peking Duck Introduction of Peking Duck: Peking Duck is a famous Chinese dish that originated in Beijing during the Imperial era. The dish dates back to the Yuan Dynasty (1271-1368) and became a staple in the Ming Dynasty (1368-1644). Traditionally, Peking Duck is known for its thin, rrispy skin and is served with pancakes, holsin sauce, and scallions. The preparation involves inflating the duck to separate the skin from the fat, marinating it, and roasting it in a closed or hung oven. It is considered a national dish of China and a symbol of Chinese culinary art. - Question: During which dynasty did the dish shown in the image become a staple in the culsine of its country? Options: (A) Tang Dynasty (B) Song Dynasty (C) Ming Dynasty (D) Qing Dynasty - Answer: (C) Ming Dynasty Image: The White House Introduction of The White House: The White House, located at 1600 Pennsylvania Avenue NW in Washington, D.C., is the official residence and workplace of the President of the United States. Construction began in 1792 and was completed in 1800. The building was designed by Trish-born architect James Hoban in the neoclassical style. It has been the residence of every U.S. president since John Adams. The White House has undergone several renovations and expansions, including the addition of the West Wing and the Oval Office. It is a symbol of the U.S. government and a site of significant historical events. - Question: Who was the architect responsible for designing the building shown in the image? Options: (A) James Hoban Now please generate the question for the Image: (concept) of (country) (introduction)</pre>	Here are two examples:					
 Question: During which dynasty did the dish shown in the image become a staple in the cuisine of its country? Options: (A) Tang Dynasty (B) Song Dynasty (C) Ming Dynasty (D) Qing Dynasty Answer: (C) Ming Dynasty Image: The White House Introduction of The White House: The White House, located at 1600 Pennsylvania Avenue NW in Washington, D.C., is the official residence and workplace of the President of the United States. Construction began in 1792 and was completed in 1800. The building was designed by Irish-born architect James Hoban in the neoclassical style. It has been the residence of every U.S. president since John Adams. The White House has undergone several renovations and expansions, including the addition of the West Wing and the Oval Office. It is a symbol of the U.S. government and a site of significant historical events. Question: Who was the architect responsible for designing the building shown in the image? Options: (A) James Hoban Now please generate the question for the Image: {concept} of {country} {introduction} 	Image: Peking Duck Introduction of Peking Duck: Peking Duck is a famous Chinese dish that originated in Beijing during the Imperial era. The dish dates back to the Yuan Dynasty (1271-1368) and became a staple in the Ming Dynasty (1368-1644). Traditionally, Peking Duck is known for its thin, crispy skin and is served with pancakes, hoisin sauce, and scallions. The preparation involves inflating the duck to separate the skin from the fat, marinating it, and roasting it in a closed or hung oven. It is considered a national dish of China and a symbol of Chinese culinary art.					
<pre>Image: The White House Introduction of The White House: The White House, located at 1600 Pennsylvania Avenue NW in Washington, D.C., is the official residence and workplace of the President of the United States. Construction began in 1792 and was completed in 1800. The building was designed by Irish-born architect James Hoban in the neoclassical style. It has been the residence of every U.S. president since John Adams. The White House has undergone several renovations and expansions, including the addition of the West Wing and the Oval Office. It is a symbol of the U.S. government and a site of significant historical events. - Question: Who was the architect responsible for designing the building shown in the image? Options: (A) James Hoban (B) Benjamin Latrobe (C) Thomas Jefferson (D) Charles Bulfinch - Answer: (A) James Hoban</pre>	 Question: During which dynasty did the dish shown in the image become a staple in the cuisine of its country? Options: (A) Tang Dynasty (B) Song Dynasty (C) Ming Dynasty (D) Qing Dynasty Answer: (C) Ming Dynasty 					
<pre>Intervalue white House Introduction of The White House: The White House, located at 1600 Pennsylvania Avenue NW in Washington, D.C., is the official residence and workplace of the President of the United States. Construction began in 1792 and was completed in 1800. The building was designed by Irish-born architect James Hoban in the neoclassical style. It has been the residence of every U.S. president since John Adams. The White House has undergone several renovations and expansions, including the addition of the West Wing and the Oval Office. It is a symbol of the U.S. government and a site of significant historical events.</pre>	Image. The White House					
- Question: Who was the architect responsible for designing the building shown in the image? Options: (A) James Hoban (B) Benjamin Latrobe (C) Thomas Jefferson (D) Charles Bulfinch - Answer: (A) James Hoban Now please generate the question for the Image: {concept} of {country} {introduction}	Introduction of The White House: The White House, located at 1600 Pennsylvania Avenue NW in Washington, D.C., is the official residence and workplace of the President of the United States. Construction began in 1792 and was completed in 1800. The building was designed by Irish-born architect James Hoban in the neoclassical style. It has been the residence of every U.S. president since John Adams. The White House has undergone several renovations and expansions, including the addition of the West Wing and the Oval Office. It is a symbol of the U.S. government and a site of significant historical events.					
Now please generate the question for the Image: {concept} of {country} {introduction}	- Question: Who was the architect responsible for designing the building shown in the image? Options: (A) James Hoban (B) Benjamin Latrobe (C) Thomas Jefferson (D) Charles Bulfinch - Answer: (A) James Hoban					
	Now please generate the question for the Image: {concept} of {country} {introduction}					

This public image shows the "{concept}" of {country}. Generate a visual reasoning multiplechoice question based on this image and the introduction of {concept}. Provide the correct answer and reason immediately following the question. 974

975 976 Here are some requirements: 977 - The question must describe a specific scenario crafted to test deeper cultural 978 understanding without directly naming {concept}. The scenario can be related to cultural 979 background, regional characteristics, historical legends, or etiquette and customs, etc. 980 - The question needs to be related to the image but does not need to describe the content of 981 the image. 982 - Ensure the question requires the user to recognize the image and use relevant knowledge to 983 answer through reasoning based on the scenario provided. Users can only answer based on the 984 image, so don't mention any "introduction" or "{concept}" in the question. 985 986 Use the following format for your introduction and question: 987 - Question: [Your Scenario-based Question] Options: (A) [Option 1] (B) [Option 2] (C) [988 Option 3] (D) [Option 4] 989 - Answer: (X) [Option X] 990 - Reason: [Your Reason for the Answer] 991 992 Here are two examples: 993 994 Image: Peking Duck 995 Introduction of Peking Duck: Peking Duck is a famous Chinese dish that originated in Beijing during the Imperial era. The dish dates back to the Yuan Dynasty (1271-1368) and became a 996 997 staple in the Ming Dynasty (1368-1644). Traditionally, Peking Duck is known for its thin, 998 crispy skin and is served with pancakes, hoisin sauce, and scallions. The preparation involves inflating the duck to separate the skin from the fat, marinating it, and roasting 999 it in a closed or hung oven. This meticulous process ensures the skin becomes crispy while 1000 1001 the meat remains tender. Peking Duck is often carved in front of diners and served in three 1002 stages: the skin, the meat, and a broth made from the bones. It is considered a national 1003 dish of {country} and a symbol of Chinese culinary art. Peking Duck has also been a part of 1004 many state banquets and diplomatic events, symbolizing Chinese hospitality and culinary 1005 excellence. 1006 1007 - Question: During a state banquet featuring the dish in the image, which aspect of its 1008 presentation is most likely emphasized to symbolize Chinese culinary excellence and 1009 hospitality? Options: (A) The use of exotic spices (B) The serving of the duck with rice (C) 1010 The incorporation of seafood (D) The carving of the duck in front of diners 1011 - Answer: (D) The carving of the duck in front of diners 1012 - Reason: The traditional carving of Peking Duck in front of diners highlights the skill 1013 involved in its preparation and serves as a symbol of Chinese culinary excellence and 1014 hospitality. 1015 1016 1017 Image: The White House 1018 Introduction of The White House: The White House, located at 1600 Pennsylvania Avenue NW in 1019 Washington, D.C., is the official residence and workplace of the President of the United 1020 States. Construction began in 1792 and was completed in 1800. The building was designed by 1021 Irish-born architect James Hoban in the neoclassical style, featuring a white-painted Aquia 1022 Creek sandstone exterior. It has been the residence of every U.S. president since John Adams 1023 . The White House has undergone several renovations and expansions, including the addition 1024 of the West Wing, East Wing, and the Oval Office. The building's iconic appearance and 1025 historical significance make it a symbol of the U.S. government and a site of significant 1026 historical events. The White House has been the location of many important decisions, 1027 meetings with foreign dignitaries, and addresses to the nation. It also serves as a museum 1028 of American history, housing numerous artifacts and pieces of art. The White House is not 1029 only a residence but also a working office, with various staff members ensuring the smooth 1030 operation of the executive branch of the U.S. government. 1031 1032 - Question: During a critical diplomatic event, the President is scheduled to meet with

CVPR 2025 Submission #00006. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

several foreign dignitaries to discuss global climate initiatives. As depicted in the image, 1033 which room inside the building is most likely to be used for this high-level diplomatic 1034 meeting? Options: (A) The Lincoln Bedroom (B) The Oval Office (C) The White House Kitchen (D 1035) The East Room 1036 - Answer: (B) The Oval Office 1037 - Reason: The Oval Office is traditionally used for important meetings and discussions, 1038 making it the most likely choice for a high-level diplomatic meeting with foreign 1039 1040 dignitaries. 1041 1042 1043 Now please generate the question for the Image: {concept} of {country}

{introduction}

6. Prompt for generating scene reasoning questions based on the introduction of cultural concepts.

	1046
(Hint: This image shows the {concept} of {country}.)	1047
	1048
Here is a question about this image:	1049
{question}	1050
	1051
First, describe the image in detail and analyze its features. Then, analyze the	1052
characteristics of the four options and compare each one with the features of the image.	1053
Finally, provide your final answer. Please include your answer into (). For example, if you	1054
choose A, please write (A).	1058
	,

7. Prompt for generating knowledge-based reasoning response to image recognition questions in training set.

(Hint: This image shows the {concept} of {country}. {introduction})	1057
	1059
Here is a question about this image:	1060
{question}	1061
	1062
First, provide a detailed description and identification of the image, analyzing its	1063
features. Then, conduct a comprehensive analysis of the question and four options based on	1064
the Hint and your knowledge. Finally, present the final answer.	1065
	1066
Please refrain from explicitly mentioning the "Hint" in your response, as these are for your	1067
discreet knowledge and not provided by the question. Please include your answer into ().	1068
For example, if you choose A, please write (A).	1069
	1070
Response:	1072

8. Prompt for generating knowledge-based reasoning response to cultural knowledge / scene reasoning questions in training set.

Here is a question about this image: Question: {question} Options: {options} First, describe and identify the image. Then, analyze the question and all four options in detail. Finally, provide the answer, indicating your final choice in parentheses. For example, if you choose A, please write (A).

9. Prompt for stepwise reasoning in the evaluation of CultureVerse.

Country	Concept	Country	Concept	Country	Concept
SUM of Concepts	19,682	Sri Lanka	90	Samoa	30
India	1,430	Democratic People's Republic of Korea	90	Turkmenistan	27
United States of America	1,411	Saudi Arabia	87	Qatar	26
Italy	661	Lithuania	86	Guyana	26
China	545	Malta	86	Kuwait	25
Mexico	524	Uzbekistan	84	Paraguay	24
Japan	522	Algeria	83	Sudan	24
Philippines	465	Lebanon	79	Angola	24
Indonesia	400	Nigeria	78	Fiji	24
France	374	Colombia	78	Seychelles	23
Russian Federation	328	Austria	77	Lesotho	22
Greece	300	Cyprus	75	Barbados	21
Germany	273	Mongolia	74	Dominica	21
Egypt	272	Cuba	73	Mauritius	20
Armenia	259	Bosnia and Herzegovina	72	Maldives	19
Australia	254	Ecuador	71	Niger	18
Spain	249	Slovakia	65	Zambia	18
Georgia	245	Iceland	65	Antigua and Barbuda	17
Brazil	239	Luxembourg	65	Saint Lucia	16
Canada	231	Iraq	62	Eswatini	16
Thailand	228	Ghana	61	Bahrain	15
Myanmar	226	Albania	60	Papua New Guinea	15
Ireland	220	Uruguay	60	Kazakhstan	15
Pakistan	218	Montenegro	60	Tuvalu	14
Nepal	216	Uganda	58	Liechtenstein	14
New Zealand	209	Chile	56	Cote d'Ivoire	14
Portugal	199	Senegal	56	Suriname	13
Ukraine	197	United Republic of Tanzania	50 50	Bahamas	12
Malaysia	189	United Arab Emirates	50 54	Entrea	12
Bangladesh	188	Guatemala	54	Niozambique Desens di	12
Peru Deland	188	Latvia	53 53	Burundi Marahali Jalanda	11
Poland	180	Afahaniatan	52	Marshall Islands	11
Bulgaria	182	Alghanistan	51	Honduras	11
Britain and Northern Ireland	176	Belarus	50	San Marino	11
Croatia	176	Banin	50	Liberia	11
Serbia	170	Oman	50	Tojikiston	11
Romania	167	Dominican Republic	30 49	Solomon Islands	10
Ethionia	155	Tunisia	48	Comoros	0
Cambodia	153	Trinidad and Tobago	40	Nauru	9
Netherlands	151	Fl Salvador	40	Vanuatu	8
Denmark	142	Monaco	47	Kiribati	8
Lao People's Democratic Republic	141	Mali	46	Timor-Leste	8
South Africa	138	Kenva	45	Grenada	8
Argentina	134	Estonia	45	Sierra Leone	7
Republic of Korea	132	Haiti	44	Rwanda	7
Czechia	130	Jamaica	43	Guinea	7
Bhutan	130	Belize	42	Libya	7
Azerbaijan	129	Plurinational State of Bolivia	41	Andorra	7
Morocco	126	Congo	40	Gambia	7
Norway	120	Namibia	39	Burkina Faso	7
Sweden	112	Somalia	39	South Sudan	5
Finland	112	Nicaragua	38	Gabon	5
Israel	109	Kyrgyzstan	36	Togo	4
Türkiye	109	Bolivarian Republic of Venezuela	35	Malawi	4
Viet Nam	109	Madagascar	35	Saint Kitts and Nevis	4
Hungary	101	Republic of Moldova	34	Djibouti	4
Ethnic_and_religiou_groups	94	Zimbabwe	33	Saint Vincent and the Grenadines	3
North Macedonia	94	Jordan	33	Central African Republic	3
Slovenia	93	Tonga	32	Chad	2
Islamic Republic of Iran	93	Botswana	32	Mauritania	2
Switzerland	93	Cameroon	31	Panama	1
Singapore	93	Costa Rica	31	Federated States of Micronesia	1
Belgium	91	Democratic Republic of the Congo	31	Equatorial Guinea	1

Table 10. Number of cultural concepts of different countries or regions