

Systematic Performance Degradation in Indic Vision-Language Models: Evidence from Hindi and Telugu

Anonymous ACL submission

Abstract

With 1.5 billion people speaking over 120 major languages, India exemplifies the challenges of multilingual AI evaluation. Current multilingual VLM benchmarks suffer from unverified auto-translations, narrow task coverage, small sample sizes, and lack of culturally grounded content. We present HinTel-AlignBench, a comprehensive evaluation framework and benchmark for Hindi and Telugu vision-language models with English-aligned samples. Our framework combines semi-automated translation with human verification to generate $\sim 4k$ QA pairs per language across five domains: adapted English datasets (VQAv2, RealWorldQA, CLEVR-Math) and native Indic sets (JEE for STEM, VAANI for cultural grounding). Evaluation of state-of-the-art open and closed-source VLMs reveals consistent performance regression from English to Indic languages, with average drops of 8.3 points for Hindi and 5.5 points for Telugu across four of five tasks. We identify key failure modes and establish reproducible baselines for multilingual multimodal evaluation.

1 Introduction

India’s 122 major languages and 1599 other languages¹ present unique challenges for multilingual AI. While recent multimodal large language models (MLLMs) such as ChatGPT (OpenAI, 2025), Gemini 2.5 (Google DeepMind, 2025), and open-weight variants (Meta Llama, 2025; Dash et al., 2025) claim multilingual support, comprehensive evaluation benchmarks for Indian languages remain scarce.

Current evaluation methodologies suffer from critical limitations in quality, scope, and scale. First, many benchmarks rely on unverified automatic translations (Wu et al., 2025), inevitably introducing noise. While text-only benchmarks

¹https://en.wikipedia.org/wiki/Languages_of_India

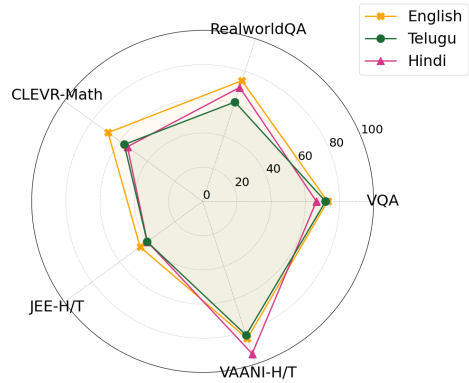


Figure 1: Average performance of GPT-4.1 and Gemini-2.5-Flash on English, Hindi, and Telugu across data-parallel visual question answering samples. Performance regresses from English to Hindi by 8.3 points and from English to Telugu by 5.5 points.

like IndicGenBench (Singh et al., 2024) exist, they lack multimodal coverage. Second, existing vision-language benchmarks often suffer from insufficient sample sizes; for instance, xChat (Yue et al., 2025) and AyaVisionBench (Dash et al., 2025) contain only 50 and 135 QA pairs per language, respectively, preventing statistically significant analysis. Third, domain coverage is often narrow. Concurrent work such as Kaleidoscope (Salazar et al., 2025) provides ~ 800 samples per Indic language but focuses exclusively on exam-based multiple-choice questions, neglecting real-world reasoning. Finally, adapted benchmarks often lack cultural grounding, evaluating surface-level translation rather than native competence (Khan et al., 2024). We elaborate on previous works in the appendix.

To address these gaps, we introduce HinTel-AlignBench, a scalable framework and benchmark for evaluating VLMs in Hindi and Telugu. Our semi-automated pipeline combines translation or LLM-based QA generation with strict human verification, achieving 5x faster processing than manual

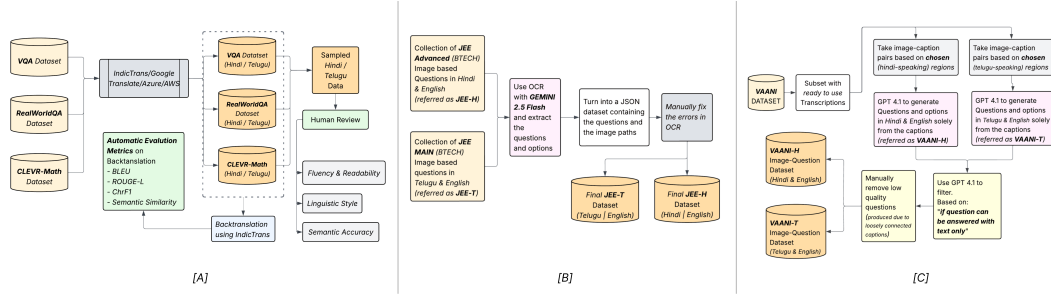


Figure 2: Dataset generation pipeline for (A) VQAv2, RealWorldQA, and CLEVR-Math using translation and human verification; (B) JEE-H and JEE-T using OCR extraction and verification; (C) VAANI-H and VAANI-T using LLM-based question generation from captions with filtering and verification.

creation for 79% of samples while maintaining linguistic fidelity. The benchmark comprises $\sim 4k$ QA pairs per language—significantly larger than prior manually verified sets—spanning five domains: real-world understanding (VQAv2 (Goyal et al., 2017a)), practical reasoning (RealWorldQA (xAI, 2024a)), visual mathematics (CLEVR-Math (Lindström and Abraham, 2022)), STEM competency (JEE-Vision from India’s Joint Entrance Exam), and cultural grounding (VAANI (Team, 2025)). Crucially, each sample includes manually verified English translations, enabling direct cross-lingual comparison.

Evaluation of state-of-the-art models on our benchmark reveals systematic performance degradation. Across all models, we observe average regressions of 8.3 points (Hindi) and 5.5 points (Telugu) relative to English, with gaps appearing in four of five tasks (Figure 1). Even frontier models like GPT-4.1 exhibit 3.8-point (Hindi) and 8.6-point (Telugu) performance drops. Performance on aligned Hindi and Telugu subsets differs by less than 1 point, indicating comparable gaps between English and both Indic languages.

Our contributions are: (1) a semi-automated framework for generating multilingual vision-language evaluation sets; (2) the largest human-verified Hindi and Telugu VLM benchmark to date, featuring culturally sourced content and English-aligned samples; and (3) a comprehensive evaluation of state-of-the-art models, highlighting significant performance regressions across diverse domains.

2 Datasets

2.1 Data Sources

We construct HinTel-AlignBench by combining translated English VQA datasets with native Indic

evaluation sets across five domains. The translated sets include 1000 samples from VQAv2 (Goyal et al., 2017b) for real-world visual understanding, 765 samples from RealWorldQA (xAI, 2024a,b) for practical reasoning, and 1000 samples from CLEVR-Math (Lindström and Abraham, 2022) for visual mathematical reasoning.

For native Indic content, we develop JEE-Vision from India’s Joint Entrance Examination, sourcing 192 Hindi questions (JEE-H) from the Advanced exam and 325 Telugu questions (JEE-T) from the Mains exam. These diagram-dependent STEM problems span mathematics, physics, and chemistry, providing the first benchmark for non-translated multilingual technical reasoning with visual content. We generate culturally grounded evaluation sets (VAANI-H and VAANI-T) by sampling 945 Hindi and 1020 Telugu images from the VAANI corpus (Team, 2025), using GPT-4.1 to create multiple-choice questions from original captions, then filtering out text-only answerable questions.

Translation-based extension enables multi-way parallel data, allowing attribution of performance to task knowledge versus language understanding (Singh et al., 2024). This approach also leverages the quality control invested in designing the original English benchmarks. Table 1 shows the distribution of QA pairs per language and task. Figure 3 showcases a few examples.

| Language | VQAv2 | RealWorldQA | CLEVR-Math | JEE-H | JEE-T | VAANI-H | VAANI-T |
|--------------|--------------|--------------|--------------|------------|------------|--------------|--------------|
| Hindi | 1,000 | 765 | 1,000 | 192 | - | 945 | - |
| Telugu | 1,000 | 765 | 1,000 | - | 325 | - | 1,020 |
| English | 1,000 | 765 | 1,000 | 192 | 325 | 945 | 1,020 |
| Total | 3,000 | 2,295 | 3,000 | 384 | 650 | 1,890 | 2,040 |

Table 1: Number of QA pairs per task per language in HinTel-AlignBench. The samples used in VQAv2, RealWorldQA and CLEVR-Math are the same across all languages.

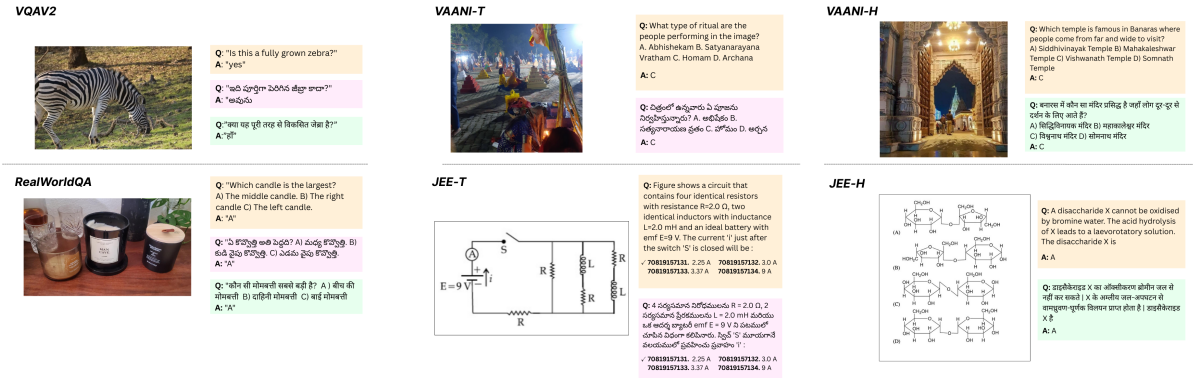


Figure 3: Qualitative Examples for different domains in our dataset. More images are shown in the appendix

2.2 Dataset Generation Framework

Figure 2 illustrates our three-stage generation framework tailored to different data sources.

Translation Pipeline. For VQAv2, RealWorldQA, and CLEVR-Math, we evaluate four translation systems (IndicTrans (Gala et al., 2023), Google Translate, Azure, AWS) on 50 diverse samples per language, selecting Azure for Hindi and AWS for Telugu. All translations undergo manual verification for semantic accuracy, linguistic style, and readability (KJ et al., 2025). We avoid back-translation-based sample selection, which introduces bias toward high-confidence translation errors. Manual review accepts 79% of VQAv2 translations without modification; among modified samples, 42% require only minor changes (verb tense), while 58% need word addition or deletion. Samples requiring only minor edits process 5x faster than generation from scratch. All verification is performed by co-author native speakers. We use one annotator per sample to maximize dataset size within budget constraints.

JEE-Vision Creation. India’s Joint Entrance Examination provides authentic STEM problems authored by subject-matter experts in target languages (JEE Mains: 13 languages; Advanced: English/Hindi), avoiding translation artifacts. We curate diagram-dependent problems, evaluating joint understanding of technical visuals and linguistic content. Questions and options are extracted using Gemini-2.5-Flash OCR (Google DeepMind, 2025), then manually verified to correct OCR errors.

VAANI Generation. From the VAANI corpus (Team, 2025), we extract images with text transcriptions from Hindi and Telugu speaking regions. Since no images have both Hindi and Tel-

ugu transcriptions, we create separate language-specific sets. Text-only GPT-4.1 generates multiple-choice questions from captions, which undergo two-stage refinement: automated filtering removes questions answerable without images, followed by human verification to eliminate low-quality questions. This process addresses cases where VAANI captions do not perfectly align with images.

3 Experimental Setup

Models. We evaluate open-weight and proprietary models claiming Indic language support. For Hindi, we test Gemma3 (4B, 12B, 27B) (Team et al., 2025), Qwen2.5VL-7B (Bai et al., 2025), Llama3.2-Vision-11B (Meta Llama, 2025), Aya-8B (Dash et al., 2025), Chittrarth-8B (Khan et al., 2024), GPT-4.1 (OpenAI, 2025), and Gemini-2.5-Flash (Google DeepMind, 2025). For Telugu, fewer models provide support; we evaluate the Gemini variants, GPT-4.1, and Chittrarth-8B.

Metrics. For multiple-choice tasks (RealWorldQA, VAANI, JEE), we report standard accuracy. For open-ended generation (VQAv2, CLEVR-Math), we utilize a hybrid evaluation protocol combining exact match with a GPT-4.1 judge to account for linguistic variations, following standard VQA practices (complete details in appendix).

4 Results and Analysis

4.1 Main Results

Tables 2 and 3 present Telugu-English and Hindi-English comparisons. Across all models and tasks, average performance regresses 5.5 points from English to Telugu and 8.3 points from English to Hindi. Performance drops occur in four of five tasks, with VAANI showing smaller gaps.

| Model | VQAv2 | | RealWorldQA | | CLEVR-Math | | JEE-T | | VAANI-T | | Ours-T | |
|------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | Tel | En | Tel | En | Tel | En | Tel | En | Tel | En | Tel | En |
| GPT 4.1 | 68.70 | 72.00 | 61.05 | 75.29 | 46.60 | <u>65.00</u> | 34.36 | 40.92 | <u>81.57</u> | 82.45 | 58.46 | 67.13 |
| Gemini 2.5 Flash | <u>75.10</u> | 74.10 | 61.18 | <u>73.20</u> | 66.70 | 71.80 | 45.90 | 54.15 | 82.75 | 80.78 | 66.33 | 70.81 |
| Gem 2.0 Flash | 70.20 | 74.20 | <u>60.92</u> | 69.67 | 43.60 | 53.50 | <u>42.15</u> | <u>53.23</u> | 80.49 | 79.61 | 59.47 | 66.04 |
| Gem 1.5 Flash | 68.50 | <u>74.40</u> | 60.00 | 67.19 | 37.40 | 46.70 | 29.85 | 39.38 | 76.27 | 79.61 | 54.40 | 61.46 |
| Chitrarth | 76.00 | 78.50 | 53.59 | 52.55 | <u>53.90</u> | 56.90 | 20.00 | 18.15 | <u>81.57</u> | <u>82.45</u> | 57.01 | 57.71 |
| Model Mean | 71.10 | 74.64 | 59.35 | 67.58 | 49.64 | 58.78 | 34.85 | 41.17 | 80.53 | 80.98 | 59.13 | 64.63 |

Table 2: Results (in %) for Telugu and English. **Bold** indicates the best and underline indicates the next best.

| Model | VQAv2 | | RealWorldQA | | CLEVR-Math | | JEE-H | | VAANI-H | | Ours-H | |
|------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | Hi | En | Hi | En | Hi | En | Hi | En | Hi | En | Hi | En |
| GPT-4.1 | 68.00 | 72.00 | 70.59 | 75.29 | 48.10 | 65.00 | <u>23.18</u> | <u>23.05</u> | <u>93.33</u> | <u>86.88</u> | 60.64 | 64.44 |
| Gemini 2.5 Flash | 65.00 | 74.10 | <u>69.54</u> | <u>73.20</u> | 60.70 | <u>71.80</u> | 56.90 | 62.89 | 93.86 | 87.19 | 69.20 | 73.84 |
| Chitrarth | <u>66.00</u> | 78.50 | 52.94 | 52.55 | <u>57.20</u> | 56.90 | 11.72 | 13.93 | 84.23 | 80.14 | 54.42 | 56.40 |
| Qwen2.5VL-7B | 37.20 | <u>74.30</u> | 51.11 | 68.10 | 29.10 | 98.80 | 17.84 | 20.70 | 81.79 | 84.76 | 43.41 | 69.33 |
| Aya-8B | 36.30 | 47.30 | 55.42 | 58.82 | 46.20 | 61.40 | 9.63 | 16.02 | 82.22 | 80.42 | 45.95 | 52.79 |
| LLaMA 3.2 11B | 35.90 | 59.80 | 35.68 | 61.57 | 18.90 | 35.60 | 14.32 | 14.45 | 77.67 | 83.28 | 36.49 | 50.94 |
| Gemma3-27B | 64.10 | 65.50 | 54.38 | 61.04 | 43.50 | 53.70 | 19.66 | 17.58 | 87.41 | 82.01 | 53.81 | 55.97 |
| Gemma3-12B | 63.00 | 65.70 | 53.98 | 58.69 | 40.20 | 46.80 | 14.84 | 16.80 | 85.50 | 82.33 | 51.50 | 54.87 |
| Gemma3-4B | 55.00 | 58.20 | 43.27 | 50.19 | 33.20 | 39.60 | 14.32 | 17.19 | 80.64 | 77.78 | 45.29 | 48.59 |
| Model Mean | 53.74 | 66.26 | 53.99 | 62.49 | 43.01 | 58.62 | 20.01 | 22.29 | 85.85 | 83.76 | 51.19 | 58.49 |

Table 3: Results (in %) for Hindi and English. **Bold** indicates the best and underline indicates the next best.

On aligned samples spanning VQAv2, CLEVR-Math, and RealWorldQA evaluated with GPT-4.1, Gemini-2.5-Flash, and Chitrarth, Hindi, Telugu, and English achieve 61.51, 62.53, and 68.6 points respectively, demonstrating systematic degradation from English to both Indic languages.

Gemini-2.5-Flash achieves best overall performance on both language pairs. Chitrarth leads on VQAv2 due to multilingual VQAv2 training. However, models show substantial cross-language variance: GPT-4.1 excels on RealWorldQA in English and Hindi but underperforms on Telugu, highlighting the need for comprehensive evaluation across all target languages. Qwen2.5VL-7B exhibits the largest Hindi-English gap at 25.92 points.

4.2 Task-Specific Analysis

Figure 4 shows average performance regression per task. CLEVR-Math and RealWorldQA exhibit the largest English-Indic gaps, while VAANI shows the smallest. VAANI-H performance exceeds English by 2.09 points on average. Analysis reveals two factors: first, some English questions fail to capture Indic-script option meanings in images with visible Indic text. Second, text-only LLM-generated distractors may enable statistical pattern exploitation.

Chain-of-Thought prompting improves reasoning tasks but benefits English more than Hindi. On JEE-H, CoT gains 13.8 points in English versus 2.22 points in Hindi, suggesting training bias to-

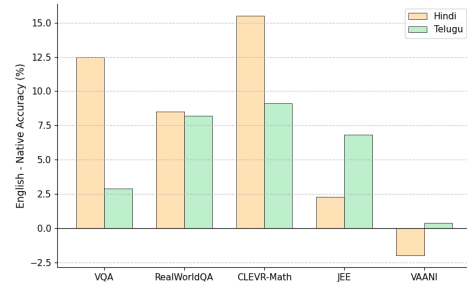


Figure 4: Average performance regression from English to Indic languages per domain. All tasks except VAANI show consistent regression.

ward English CoT data. Detailed ablation studies and error analysis are reported in the appendix.

5 Conclusion

This paper introduces HinTel-AlignBench, a framework for developing benchmarks to evaluate multimodal large language models in Hindi and Telugu, addressing critical gaps in existing multilingual evaluations. We combined semi-automated dataset creation with rigorous human verification and sourced culturally grounded native datasets to assess diverse capabilities. Evaluations of state-of-the-art VLMs reveal significant performance regressions in Indic languages compared to English, emphasizing the need for targeted improvements in multilingual visual understanding.

245 Limitations

246 While our benchmark introduces a diverse evaluation
247 set it has limitations. First, the proprietary
248 models we evaluated achieve high scores on the
249 VAANI-H/T. We use text only LLMs for generat-
250 ing QA from VAANI captions and they often do
251 not design good distractors. Thus, the models may
252 guess the correct answer by exploiting statistical
253 patterns, which inflates metrics. A future work is
254 using Multi-Binary Accuracy (Cai et al., 2024) for
255 VAANI subsets. Second, the JEE-H benchmark
256 contains only 2 Mathematics questions, due to lack
257 of image-based mathematics questions in the JEE-
258 Advanced examination. Finally, there are 22 of-
259 ficial Indian languages and we cover English and
260 2/22 (Hindi and Telugu) with this work. We hope
261 this benchmark gets extended to all other Indic
262 languages with contributions from native speakers
263 from those languages.

264 References

265 Nahid Alam, Karthik Reddy Kanjula, Surya Guthikonda,
266 Timothy Chung, Bala Krishna S Vegesna, Abhipsha
267 Das, Anthony Susevski, Ryan Sze-Yin Chan, SM Ud-
268 din, Shayekh Bin Islam, and 1 others. 2025. Behind
269 maya: Building a multilingual vision language model.
270 *arXiv preprint arXiv:2505.08910*.

271 Anthropic. 2025. Claude 3.7 sonnet:
272 The first hybrid reasoning model. An-
273 thropic Company Website. Available at:
274 [https://www.anthropic.com/news/claude-3-7-](https://www.anthropic.com/news/claude-3-7-sonnet)
275 [sonnet](https://www.anthropic.com/news/claude-3-7-sonnet).

276 Daman Arora, Himanshu Singh, and Mausam. 2023.
277 [Have LLMs advanced enough? a challenging prob-](#)
278 [lem solving benchmark for large language models.](#)
279 In *Proceedings of the 2023 Conference on Empiri-*
280 *cal Methods in Natural Language Processing*, pages
281 7527–7543, Singapore. Association for Computa-
282 tional Linguistics.

283 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wen-
284 bin Ge, Sibao Song, Kai Dang, Peng Wang, Shi-
285 jie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu,
286 Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei
287 Wang, Wei Ding, Zheren Fu, Yiheng Xu, and 8 others.
288 2025. [Qwen2.5-vl technical report](#). *arXiv preprint*
289 *arXiv:2502.13923*.

290 Mu Cai, Reuben Tan, Jianrui Zhang, Bocheng Zou, Kai
291 Zhang, Feng Yao, Fangrui Zhu, Jing Gu, Yiwu Zhong,
292 Yuzhang Shang, Yao Dou, Jaden Park, Jianfeng Gao,
293 Yong Jae Lee, and Jianwei Yang. 2024. Temporal-
294 bench: Towards fine-grained temporal understand-
295 ing for multimodal video models. *arXiv preprint*
296 *arXiv:2410.10818*.

Soravit Changpinyo, Linting Xue, Michal Yarom, 297
Ashish V Thapliyal, Idan Szpektor, Julien Amelot, 298
Xi Chen, and Radu Soricut. 2023. [MaXM: Towards](#)
299 [multilingual visual question answering](#). In *The 2023*
300 *Conference on Empirical Methods in Natural Lan-*
301 *guage Processing*. 302

Wenliang Dai, Junnan Li, Dongxu Li, Anthony Tiong, 303
Junqi Zhao, Weisheng Wang, Boyang Li, Pascale 304
Fung, and Steven Hoi. 2023. [InstructBLIP: Towards](#)
305 [general-purpose vision-language models with instruc-](#)
306 [tion tuning](#). In *Thirty-seventh Conference on Neural*
307 *Information Processing Systems*. 308

Saurabh Dash, Yiyang Nan, John Dang, Arash Ah- 309
madian, Shivalika Singh, Madeline Smith, Bharat 310
Venkitesh, Vlad Shmyhlo, Viraat Aryabumi, Walter 311
Beller-Morales, Jeremy Pektmez, Jason Ozuzu, Pierre 312
Richemond, Acyr Locatelli, Nick Frosst, Phil Blun- 313
som, Aidan Gomez, Ivan Zhang, Marzieh Fadaee, 314
and 6 others. 2025. [Aya vision: Advancing the fron-](#)
315 [tier of multilingual multimodality](#). *arXiv preprint*
316 *arXiv:2505.08751*, arXiv:2505.08751. 317

Matt Deitke, Christopher Clark, Sangho Lee, Rohun 318
Tripathi, Yue Yang, Jae Sung Park, Mohammadreza 319
Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, 320
Jiasen Lu, Taira Anderson, Erin Bransom, Kiana 321
Ehsani, Huong Ngo, YenSung Chen, Ajay Patel, 322
Mark Yatskar, Chris Callison-Burch, and 31 oth- 323
ers. 2025. [Molmo and pixmo: Open weights and](#)
324 [open data for state-of-the-art vision-language mod-](#)
325 [els](#). *Conference on Computer Vision and Pattern*
326 *Recognition (CVPR)*. 327

Jay Gala, Pranjal A Chitale, A K Raghavan, Varun 328
Gumma, Sumanth Doddapaneni, Aswanth Kumar M, 329
Janki Atul Nawale, Anupama Sujatha, Ratish Pudup- 330
pully, Vivek Raghavan, Pratyush Kumar, Mitesh M 331
Khapra, Raj Dabre, and Anoop Kunchukuttan. 2023. 332
[Indictrans2: Towards high-quality and accessible ma-](#)
333 [chine translation models for all 22 scheduled indian](#)
334 [languages](#). *Transactions on Machine Learning Re-*
335 *search*. 336

Google DeepMind. 2025. [Gemini 2.5: Our most intelli-](#)
337 [gent ai model](#). Google DeepMind Blog. Released on
338 March 25, 2025. 339

Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv 340
Batra, and Devi Parikh. 2017a. Making the V in VQA 341
matter: Elevating the role of image understanding 342
in Visual Question Answering. In *Conference on*
343 *Computer Vision and Pattern Recognition (CVPR)*. 344

Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv 345
Batra, and Devi Parikh. 2017b. Making the v in vqa 346
matter: Elevating the role of image understanding 347
in visual question answering. In *Proceedings of the*
348 *IEEE Conference on Computer Vision and Pattern*
349 *Recognition (CVPR)*, pages 6904–6913. 350

Drew A Hudson and Christopher D Manning. 2019. 351
Gqa: A new dataset for real-world visual reason- 352
ing and compositional question answering. *Confer-*
353

| | | | | | |
|-----|--|--|--|--|-----|
| 354 | | | | | |
| 355 | | | | | |
| 356 | Shaharukh Khan, Ayush Tarun, Abhinav Ravi, Ali Faraz, | | | | |
| 357 | Akshat Patidar, Praveen Kumar Pokala, Anagha | | | | |
| 358 | Bhangare, Raja Kolla, Chandra Khatri, and Shubham | | | | |
| 359 | Agarwal. 2024. Chitrarth: Bridging vision and lan- | | | | |
| 360 | guage for a billion people. In <i>NeurIPS Multimodal</i> | | | | |
| 361 | <i>Algorithmic Reasoning</i> . | | | | |
| 362 | Sankalp KJ, Ashutosh Kumar, Laxmaan Balaji, Nikunj | | | | |
| 363 | Kotecha, Vinija Jain, Aman Chadha, and Sreyoshi | | | | |
| 364 | Bhaduri. 2025. Indicmmlu-pro: Benchmarking ind- | | | | |
| 365 | ic large language models on multi-task language | | | | |
| 366 | understanding. <i>arXiv preprint arXiv:2501.15747</i> . | | | | |
| 367 | Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. | | | | |
| 368 | 2023. Blip-2: Bootstrapping language-image pre- | | | | |
| 369 | training with frozen image encoders and large lan- | | | | |
| 370 | guage models. In <i>International conference on ma-</i> | | | | |
| 371 | <i>chine learning</i> , pages 19730–19742. PMLR. | | | | |
| 372 | Adam Dahlgren Lindström and Savitha Sam Abraham. | | | | |
| 373 | 2022. Clevr-math: A dataset for compositional lan- | | | | |
| 374 | guage, visual, and mathematical reasoning. <i>arXiv</i> | | | | |
| 375 | <i>preprint</i> . | | | | |
| 376 | Fangyu Liu, Emanuele Bugliarello, Edoardo Maria | | | | |
| 377 | Ponti, Siva Reddy, Nigel Collier, and Desmond El- | | | | |
| 378 | liott. 2021. Visually grounded reasoning across lan- | | | | |
| 379 | guages and cultures. In <i>Proceedings of the 2021</i> | | | | |
| 380 | <i>Conference on Empirical Methods in Natural Lan-</i> | | | | |
| 381 | <i>guage Processing</i> , pages 10467–10485, Online and | | | | |
| 382 | Punta Cana, Dominican Republic. Association for | | | | |
| 383 | Computational Linguistics. | | | | |
| 384 | Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae | | | | |
| 385 | Lee. 2023. Visual instruction tuning. <i>Advances in</i> | | | | |
| 386 | <i>neural information processing systems</i> , 36:34892– | | | | |
| 387 | 34916. | | | | |
| 388 | Meta Llama. 2025. Llama 3.2-vision: Instruction-tuned | | | | |
| 389 | image reasoning generative models. Model release | | | | |
| 390 | by Meta. Available at: https://huggingface.co/meta- | | | | |
| 391 | llama/Llama-3.2-11B-Vision . | | | | |
| 392 | OpenAI. 2025. Gpt-4.1: A new series of gpt models | | | | |
| 393 | with major improvements on coding, instruction fol- | | | | |
| 394 | lowing, and long context. OpenAI Company Website. | | | | |
| 395 | Released on April 14, 2025. | | | | |
| 396 | Jonas Pfeiffer, Gregor Geigle, Aishwarya Kamath, Jan- | | | | |
| 397 | Martin O. Steitz, Stefan Roth, Ivan Vulić, and Iryna | | | | |
| 398 | Gurevych. 2022. xGQA: Cross-lingual visual ques- | | | | |
| 399 | tion answering. In <i>Findings of the Association for</i> | | | | |
| 400 | <i>Computational Linguistics: ACL 2022</i> , pages 2497– | | | | |
| 401 | 2511, Dublin, Ireland. Association for Computational | | | | |
| 402 | Linguistics. | | | | |
| 403 | Hanoona Rasheed, Muhammad Maaz, Abdelrahman | | | | |
| 404 | Shaker, Salman Khan, Hisham Cholakkal, Rao M. An- | | | | |
| 405 | wer, Tim Baldwin, Michael Felsberg, and Fahad S. | | | | |
| 406 | Khan. 2025. Palo: A polyglot large multimodal | | | | |
| 407 | model for 5b people. In <i>Winter Conference on Ap-</i> | | | | |
| 408 | <i>plications of Computer Vision (WACV)</i> , pages 1745– | | | | |
| 409 | 1754. | | | | |
| | David Romero, Chenyang Lyu, Haryo Akbarianto Wi- | | | | 410 |
| | bowo, Teresa Lynn, Injy Hamed, Aditya Nanda | | | | 411 |
| | Kishore, Aishik Mandal, Alina Dragonetti, Artem | | | | 412 |
| | Abzaliev, Atnafu Lambebo Tonja, Bontu Fufa Balcha, | | | | 413 |
| | Chenxi Whitehouse, Christian Salamea, Dan John | | | | 414 |
| | Velasco, David Ifeoluwa Adelani, David Le Meur, | | | | 415 |
| | Emilio Villa-Cueva, Fajri Koto, Fauzan Farooqui, | | | | 416 |
| | and 56 others. 2024. Cvqa: Culturally-diverse mul- | | | | 417 |
| | tilingual visual question answering benchmark. In | | | | 418 |
| | <i>Advances in Neural Information Processing Systems</i> , | | | | 419 |
| | volume 37, pages 11479–11505. Curran Associates, | | | | 420 |
| | Inc. | | | | 421 |
| | Israfil Salazar, Manuel Fernández Burda, Shayekh Bin | | | | 422 |
| | Islam, Arshia Soltani Moakhar, Shivalika Singh, | | | | 423 |
| | Fabian Farestam, Angelika Romanou, Danylo | | | | 424 |
| | Boiko, Dipika Khullar, Mike Zhang, Dominik | | | | 425 |
| | Krzemiński, Jekaterina Novikova, Luísa Shimabu- | | | | 426 |
| | coro, Joseph Marvin Imperial, Rishabh Maheshwary, | | | | 427 |
| | Sharad Duwal, Alfonso Amayuelas, Swati Rajwal, | | | | 428 |
| | Jebish Purbey, and 25 others. 2025. Kaleidoscope: | | | | 429 |
| | In-language exams for massively multilingual vision | | | | 430 |
| | evaluation. <i>Preprint</i> , arXiv:2504.07072. | | | | 431 |
| | Harman Singh, Nitish Gupta, Shikhar Bharadwaj, Di- | | | | 432 |
| | nesh Tewari, and Partha Talukdar. 2024. IndicGen- | | | | 433 |
| | Bench: A multilingual benchmark to evaluate gen- | | | | 434 |
| | eration capabilities of LLMs on Indic languages. In | | | | 435 |
| | <i>Proceedings of the 62nd Annual Meeting of the As-</i> | | | | 436 |
| | <i>sociation for Computational Linguistics (Volume 1:</i> | | | | 437 |
| | <i>Long Papers)</i> , pages 11047–11073, Bangkok, Thai- | | | | 438 |
| | land. Association for Computational Linguistics. | | | | 439 |
| | Jingqun Tang, Qi Liu, Yongjie Ye, Jinghui Lu, Shu | | | | 440 |
| | Wei, Chunhui Lin, Wanqing Li, Mohamad Fitri | | | | 441 |
| | Faiz Bin Mahmood, Hao Feng, Zhen Zhao, and 1 oth- | | | | 442 |
| | ers. 2024. Mtvqa: Benchmarking multilingual text- | | | | 443 |
| | centric visual question answering. <i>arXiv preprint</i> | | | | 444 |
| | <i>arXiv:2405.11985</i> . | | | | 445 |
| | Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya | | | | 446 |
| | Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, | | | | 447 |
| | Tatiana Matejovicova, Alexandre Ramé, Morgane | | | | 448 |
| | Rivière, and 1 others. 2025. Gemma 3 technical | | | | 449 |
| | report. <i>arXiv preprint arXiv:2503.19786</i> . | | | | 450 |
| | VAANI Team. 2025. Vaani: Capturing the language | | | | 451 |
| | landscape for an inclusive digital india (phase 1). | | | | 452 |
| | https://vaani.iisc.ac.in/ . | | | | 453 |
| | Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten | | | | 454 |
| | Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, | | | | 455 |
| | and 1 others. 2022. Chain-of-thought prompting elic- | | | | 456 |
| | its reasoning in large language models. <i>Advances</i> | | | | 457 |
| | <i>in neural information processing systems</i> , 35:24824– | | | | 458 |
| | 24837. | | | | 459 |
| | Minghao Wu, Weixuan Wang, Sinuo Liu, Huifeng | | | | 460 |
| | Yin, Xintong Wang, Yu Zhao, Chenyang Lyu, | | | | 461 |
| | Longyue Wang, Weihua Luo, and Kaifu Zhang. | | | | 462 |
| | 2025. The bitter lesson learned from 2,000+ multilin- | | | | 463 |
| | gual benchmarks. <i>arXiv preprint arXiv:2504.15521</i> , | | | | 464 |
| | arXiv:2504.15521. | | | | 465 |
| | xAI. 2024a. Grok-1.5 vision preview. https://x.ai/ | | | | 466 |
| | blog/grok-1.5v . | | | | 467 |

468 xAI. 2024b. [Realworldqa dataset](#). Hugging Face
469 Dataset Repository.

470 XAI.org, 2024. Realworldqa: A benchmark for real-
471 world spatial understanding capabilities of mul-
472 timodal ai models. [https://huggingface.co/
473 datasets/xai-org/RealWorldQA](https://huggingface.co/datasets/xai-org/RealWorldQA). RealWorldQA
474 is a benchmark designed for real-world understand-
475 ing with 765 multiple-choice questions requiring
476 recognition of details in high-resolution images.

477 Xiang Yue, Yueqi Song, Akari Asai, Seungone Kim,
478 Jean de Dieu Nyandwi, Simran Khanuja, Anjali Kan-
479 tharuban, Lintang Sutawika, Sathyanarayanan Ra-
480 mamoorthy, and Graham Neubig. 2025. [Pangea: A
481 fully open multilingual multimodal LLM for 39 lan-
482 guages](#). In *The Thirteenth International Conference
483 on Learning Representations*.

484 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
485 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,
486 Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang,
487 Joseph E. Gonzalez, and Ion Stoica. 2023. Judging
488 llm-as-a-judge with mt-bench and chatbot arena. In
489 *Advances in neural information processing systems*.

490 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and
491 Mohamed Elhoseiny. 2024. [MiniGPT-4: Enhancing
492 vision-language understanding with advanced large
493 language models](#). In *The Twelfth International Con-
494 ference on Learning Representations*.

495 Jinguo Zhu and 1 others. 2025. [Internvl3: Explor-
496 ing advanced training and test-time recipes for
497 open-source multimodal models](#). *arXiv preprint
498 arXiv:2504.10479*.

A Appendix

A.1 Qualitative Examples

Refer to Fig. 5, with examples from each of the datasets, i.e VQAv2, RealWorldQA, CLEVR-Math, VAANI-H, JEE-T, VAANI-T, JEE-H

A.2 Evaluation Metrics

The RealWorldQA, VAANI-H/T subsets have only multiple-choice questions and JEE-T has multiple-choice or integer-answer questions. We use accuracy as the evaluation metric for all of these sets. We extract the answers using regex-based parsing, and report the overall accuracy across all the questions.

Hybrid Evaluation for VQA and CLEVR-Math: For VQAv2 and CLEVR-Math subsets, the answers are either a single word or short phrases. We adopt a hybrid evaluation strategy. We first evaluate a sample using exact match. Our exact match evaluation is built using the official VQA evaluation script (Goyal et al., 2017a), with the functionalities also extended to Hindi and Telugu. While exact match is strict and interpretable, it may penalize correct answers with minor surface-level variations (e.g., “yes” and “yes, it is”, synonyms, etc.). If exact match fails for a sample, we evaluate that sample using "gpt-4.1-2025-04-14" (OpenAI, 2025) as described in (Zheng et al., 2023). This two-step approach enables both high precision and flexibility, especially in cases in which answers may vary in form but not meaning. Its impact on scores is discussed in the appendix.

JEE-H Evaluation: The JEE-H Dataset has single correct MCQs, multiple correct MCQs, and numeric-type questions. We extend the scoring process used in (Arora et al., 2023). However, instead of the manual step they use, we replace with regex-based parsing to extract answers, and rule based processing to score each question. See appendix for details.

A.3 Ablation Study

We ablate prompting and other baselines using the Gemma-27B model on the Hindi subsets. We report comparisons for Direct Inference (Standard Prompting), Chain-of-Thought (Wei et al., 2022), and a Caption-Only (Deitke et al., 2025) baseline. Table 4 contains the results of ablation study.

Standard prompting provides a strong baseline performance across all datasets, especially on VAANI and VQA, indicating that the model

benefits significantly from both modalities in relatively straightforward recognition tasks. In contrast, Chain-of-Thought (CoT) prompting notably improves performance on reasoning-heavy datasets. These improvements highlight the effectiveness of step-by-step reasoning prompts in tasks that require deeper cognitive processing. However, CoT does not benefit all tasks, with marginal gains or regression on VAANI and VQAv2, where reasoning depth is less critical. Overall, CoT benefits the English subset by 3.78 points in comparison to the benefit on Hindi at 0.79 points. This indicates a possible bias in the training data, with possibly more training done on English CoT data.

Next, to explore the need for visual embedding for QA, we evaluate a caption-only baseline. Here, the Gemma-27B model is asked to first generate a rich textual summary of the image. Then the generated caption and the original question are given to the model without the original image. This approach is inferior to the direct image-based inputs in all tasks.

A.4 Failure Analysis of Visual Language Models

We categorized the failures of GPT 4.1 on VAANI-T subset for both English and Telugu. We analyze results on VAANI-T as it sources images and QA from the Indian context and failures on this set point to common reasons for errors for current multimodal multilingual LLMs. We found 4 major categories of errors that appeared in both English and Telugu questions. They are: (1) Missing Indian context: Failed predictions in this category were due to missing knowledge specific to India needed to answer the question. (2) Visual Grounding error: Failures in this category were due to the model’s inability to associate objects in the picture with the question or options. (3) Visual Perception Failure: Errors in this category are due to the model failing to interpret the content of the image correctly. (4) Failure to Ground in Indian Context: Samples in this category accurately identified visual elements but they failed to analyze them in the appropriate socio-cultural context. Table 6 reports the percentages of the different categories on this dataset.

A.5 Related Work

Multimodal Instruction-Tuned Models. Previous work has shown that large language models can be extended to vision tasks by instruction tuning with multimodal data (Liu et al., 2023; Dai et al.,

| Method | VQAv2 | | RealWorldQA | | JEE-H | | CLEVR-Math | | VAANI-H | | Ours-H | |
|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | Hi | En | Hi | En | Hi | En | Hi | En | Hi | En | Hi | En |
| Standard | 64.10 | 65.50 | <u>54.38</u> | <u>61.04</u> | <u>19.66</u> | <u>17.58</u> | <u>43.50</u> | <u>53.70</u> | 87.41 | <u>82.01</u> | 53.81 | 55.97 |
| CoT | <u>60.50</u> | 62.70 | 61.96 | 64.31 | 21.88 | 31.38 | 44.3 | 57.1 | <u>84.34</u> | 83.28 | 54.60 | 59.75 |
| Caption-Only | 56.40 | <u>64.1</u> | 44.05 | 54.77 | – | – | 36.4 | 32.9 | 80.85 | 77.14 | – | – |

Table 4: Results on Gemma-27B across multiple datasets using different inference methods. **Bold** indicates the best performance, underline indicates the second best.

| Benchmark | Indic Lang | QA Type | Image count | QA | Indic QA per lang | Human | Culturally Sourced |
|----------------|---|---------|-------------|-------|-------------------|-------|--------------------|
| AyaVisionBench | hin | Chat | 3.1k | 3.1k | 135 | ✗ | ✗ |
| xMMMU | hin | MC | 300 | 3k | 291 | ✗ | ✗ |
| xGQA | ben | OE | 300 | 77.3k | 9.7k | ✗ | ✗ |
| MTVQA | - | OE | 8.8k | 28.6k | - | ✓ | ✗ |
| M3Exam | - | MC | 2.8k | 12.3k | - | ✓ | ✗ |
| xChatBench | hin | Chat | 400 | 400 | 50 | ✓ | ✗ |
| Kaleidoscope | ben, hin , tel | MC | 20.9k | 20.9k | ~ 800 | ✓ | ✗ |
| XM100 | ben, hin , tel | Caption | 100 | 3.6k | 100 | ✓ | ✓ |
| MaRVL | tam | Caption | 5.5k | 5.7k | 1.2k | ✓ | ✓ |
| MaXM | hin | OE | 1.4k | 2.1k | 294 | ✓ | ✓ |
| CVQA | ben, urd, tam, mar, hin , tel | MC | 5.2k | 10.4k | ~ 300 | ✓ | ✓ |
| Ours | hin , tel | OE, MC | 5.1k | 13.5k | ~ 4k | ✓ | ✓ |

Table 5: Comparison of existing multilingual Visual Question Answering (VQA) benchmarks with ours. “QA Type” denotes the question-answering format (Chat = Multimodal chat; OE = open-ended VQA; MC = multiple-choice VQA; Caption = captioning/multi-image captioning); “QA” denotes the total question-answer pairs; “Indic QA per lang” denotes the question-answer pairs per Indic language; “Human” indicates human verified/annotated data; “Culturally Sourced” indicates culturally-sourced data. xChatBench, xMMMU and XM100 are part of PangeaBench (Yue et al., 2025). Ours (bottom row) is the only VQA dataset with human-verified annotations, diverse culturally sourced samples and a significant sample size per Indic language.

| Error Category | English | Telugu |
|-------------------------------------|---------|--------|
| Lack of Knowledge about India | 17% | 16% |
| Visual Grounding Error | 49% | 50% |
| Visual Perception Failure | 19% | 18% |
| Failure to Ground in Indian Culture | 15% | 16% |

Table 6: Error category distribution for GPT 4.1 on VAANI-T for English and Telugu. Percentages indicate the proportion of errors falling into each category.

2023; Li et al., 2023; Zhu et al., 2024). Many open weight Vision-Language Models such as Gemma3 (Team et al., 2025), Qwen2.5VL (Bai et al., 2025), InternVL3 (Zhu et al., 2025), Molmo (Deitke et al., 2025), Llama3.2 Vision (Meta Llama, 2025) and closed source models such as GPT-4.1 (OpenAI, 2025), Claude 3.7 (Anthropic, 2025), Gemini 2.5 (Google DeepMind, 2025)) have since been released, all leveraging some form of multimodal instruction tuning. Many of these recent models are multilingual, but the evaluation of their multilingual abilities is limited and is not reported on

a consistent and reliable multimodal multilingual benchmark (Dash et al., 2025; Yue et al., 2025; Rasheed et al., 2025; Alam et al., 2025).

Multilingual Vision-Language Benchmarks.

There have been a number of recent benchmarks which evaluate multi-lingual abilities and cultural robustness of VLMs. CVQA (Romero et al., 2024) contains 10.4k visual QA pairs across 31 languages, demonstrating large performance gaps on low-resource languages. MTVQA (Tang et al., 2024) provides text-centric VQA in 9 languages with human annotations. xGQA (Pfeiffer et al., 2022) automatically translates the GQA (Hudson and Manning, 2019) dataset into 7 languages while MaxM (Changpinyo et al., 2023) uses Machine Translation plus lightweight post-editing for 7 languages, with less than 300 samples per language. MarVL (Liu et al., 2021) focuses on culturally sourced reasoning by evaluating whether a caption about an image pair is true or false, but is not a Visual Question Answering dataset. Concurrently with our work, Kaleidoscope (Salazar et al., 2025) is a multi-

lingual multimodal benchmark with 800 VQA MCQs per Indic language sourced from regional examinations. However, Kaleidoscope has a narrow scope as it focuses only on exam-based questions and does not evaluate real-world understanding and practical reasoning. Multimodal benchmarks released along with their corresponding models include PangeaBench (Yue et al., 2025) (14 datasets in 47 languages), and AyaVisionBench (Dash et al., 2025) (9 tasks in 23 languages). However, these sets commonly have few manually verified samples per language, making the results less reliable for a particular target language. Table 5 provides an overview and compares them with our benchmark. A significant portion of our images per language (>1K) and our QA (>1K) are using culturally sourced, which is 5 to 20x larger than the culturally relevant subsets of previous datasets in this field.

Indic and Domain-Specific Datasets. Although there are dozens of English VQA resources (VQAv2 (Goyal et al., 2017a), CLEVR-Math (Lindström and Abraham, 2022), RealWorldQA (XAI.org, 2024)), there is a need for benchmarks in Indian languages. Previous works include (Singh et al., 2024; Arora et al., 2023). IndicGenBench (Singh et al., 2024) consists of 5 diverse tasks in 29 Indic languages but does not include multimodality. JEEBench (Arora et al., 2023) consists of questions from the highly competitive IIT JEE-Advanced exam, but does not include images-related questions. In contrast, our JEEset only includes questions with images. Our work fills the need for an accurate and diverse multimodal benchmark by providing ~4k Visual Question Answering each in Hindi and Telugu.

A.6 Compute and Costs

We primarily ran all open-source models on H100 and A100 GPUs, rented via the Akash Console Network², incurring approximately 100 USD in compute costs. For proprietary models, we used APIs, leveraging Gemini’s free tier and spending around 60 USD on OpenAI’s APIs. The total expenditure amounts to approximately 160 USD.

A.7 Prompts

JEE-H COT Inference

`eng_prompt = ""`"You are an expert at solving physics, chemistry and math questions that involve both text and images. Analyze the text and the associated image to answer the question above.

²<https://akash.network/>

First, think step-by-step and provide your detailed reasoning. Explain how you interpret the text and the image, the principles or formulas you are using, and the intermediate calculations or logical steps you take to arrive at the solution.

After your reasoning, provide the final answer on a new line.

Your entire response, including the reasoning and the final answer, MUST end with the following line exactly:

`{ answer: }""`

JEE-H/T Standard Inference

`eng_prompt = "You are an expert at solving physics, chemistry and math questions that involve both text and images.`

Analyze the text and the associated image to answer the question above. Provide only the final answer, without any explanations or intermediate steps. Your response MUST end with the following line: `{ answer: }`"

VQA LLM Eval

`System_Prompt = "You are an expert judge evaluating semantic similarity between two answers.`

Respond only with '1' for similar or '0' for not similar."

`User_Prompt = "Are the following two answers semantically similar?"`

Answer 1: `{a}`

Answer 2: `{b}`

Respond with only '1' if they are similar in meaning, or '0' if they are not similar."

CLEVR LLM Eval

`System_Prompt = "You are an expert judge evaluating semantic similarity between two answers.`

Respond only with '1' if they are semantically similar, or '0' if they are not."

`User_Prompt = "Are the following two {lang} answers semantically similar?"`

Answer 1: `{ans1}`

Answer 2: `{ans2}`

Ignore minor differences in phrasing. Respond with '1' if they express the same meaning, or '0' otherwise."

A.8 Inference Strategies

This appendix details the various inference techniques employed for evaluating the multimodal models, along with the specific prompt setups used for English and Hindi. Below are setups in English.

Caption-Only Baseline

This technique involves a two-step process:

1) Dense Caption Generation: The model first

688 generates a detailed, factual description (dense cap-
689 tion) of the provided image.

690 **2) Question Answering with Caption:** The model
691 then answers the original question using *only* the
692 generated dense caption as context, without access
693 to the original image.

694 Prompt for Step 1 (Dense Caption Generation):

eng_prompt = Don't forget these rules:

1. Be Direct and Concise: Provide straightforward descriptions without adding interpretative or speculative elements.

2. Use Segmented Details: Break down details about different elements of an image into distinct sentences, focusing on one aspect at a time.

3. Maintain a Descriptive Focus: Prioritize purely visible elements of the image, avoiding conclusions or inferences.

4. Follow a Logical Structure: Begin with the central figure or subject and expand outward, detailing its appearance before addressing the surrounding setting.

5. Avoid Juxtaposition: Do not use comparison or contrast language; keep the description purely factual.

6. Incorporate Specificity: Mention age, gender, race, and specific brands or notable features when present, and clearly identify the medium if it's discernible.

When writing descriptions, prioritize clarity and direct observation over embellishment or interpretation. Write a detailed description of this image, do not forget about the texts on it if they exist. Also, do not forget to mention the type/style of the image. No bullet points. Start with the words, "This image displays:"

695 Prompt for Step 2 (Question Answering with
696 Caption):

Context: {generated_caption}

Question: {question_text}

Answer:

697 **Chain of Thought (CoT) Prompting**

698 This technique guides the model to generate a step-
699 by-step thought process (rationale) before arriving
700 at the final answer. This is intended to improve
701 reasoning capabilities for complex questions.

702 The image is provided along with language-
703 specific system and user prompts.

English System Prompt:

When provided with an image and a question, generate a rationale first and then derive an answer.

Your rationale should include detailed visual elements in order to derive the answer.

English User Prompt:

Answer the question with following instruction:

1. Generate a rationale first and then derive an answer.

2. For your final answer, provide a Correct Option Number Only.

Question:

{question}

Output Format

<rationale>

Answer: <your answer>

Standard (Direct Inference)

704

The model is provided with both the image and the question and is expected to generate a direct answer without explicit intermediate reasoning steps requested in the prompt.

705

706

707

708

The image and the question are provided to the model. The prompt typically instructs the model to answer directly, often in a specific format (e.g., option number for multiple-choice questions).

709

710

711

712

Example input (for a multiple-choice question):

713

Image: [Image Data]

Question: {question_text_with_options}

Instruction: Provide the Correct Option Number Only.

Answer:

(Note: This often shares the final answer formatting instruction with CoT, but without the explicit rationale generation step.)

714

715

716

JEE-H Scoring Rules

717

- For single correct Multi-Choice Questions (MCQs) and integer-type questions, a binary score was assigned: 1 if the model answered with the correct option/integer and 0 otherwise.

718

719

720

721

722

- For multi-correct MCQs, a score of 1 is given when the generated response consists of all and only the correct options. If any of the wrong option is chosen, the response is given a score of 0. If the response contains no incorrect option and some of the correct options, a score of 0.25 is given for each of the correct options.

723

724

725

726

727

728

729

730

- For numeric-type questions, answers in the range of ± 0.01 with the gold answer are given a score of 1, and 0 otherwise.

731

732

733

- In Chain-of-Thought (CoT) inferencing, if a Hindi response included reasoning in English, 0 was given irrespective of the answer.

734

735

736

- In standard zero-shot inference, where the model was explicitly instructed to return only a direct answer, responses containing any explanation or reasoning were also scored 0, irrespective of the answer.

737

738

739

740

741

Impact of LLM Judge on Evaluation Scores

As detailed in Section A.2, our hybrid evaluation strategy for CLEVR-Math and VQA incorporates an LLM judge to re-evaluate answers initially marked incorrect by the exact match protocol. This approach revealed a notable trend: the score increase attributed to the LLM judge was considerably more for the responses in Hindi and Telugu as compared to English. For instance, when evaluating the VQA dataset, the average scores for English responses increased by $\approx 4\%$ after LLM judging. In contrast, Hindi scores on the same dataset increased by $\approx 12\%$ after the initial exact match scores. This disparity suggests that the models demonstrate a greater proficiency in adhering to strict answer formatting conventions in English, while responses in Hindi and Telugu, though often semantically correct, are more frequently penalized by the exact match criteria due to variations in phrasing or formatting. The LLM judge effectively mitigates this by recognizing a broader range of correct expressions.

A.9 Licenses

The datasets used in our benchmark are released under the following licenses:

- VQAv2, VAANI and ClevrMath are released under Creative Commons Attribution 4.0 International (CC BY 4.0) license.
- RealWorldQA is released under Creative Commons Attribution No Derivatives 4.0.
- JEE-Vision is built from publicly available JEE questions. We will be releasing links under CC BY-NC 4.0 license along with the QA pair. The exam papers are distributed for the general public (<https://www.jeeadv.ac.in/archive.html>) used by various publishers and institutes for commercial purposes already. When collecting this data, we followed the steps by previous works on JEE Bench (Arora et al., 2023).

