# Kaleidoscope: In-language Exams for Massively Multilingual Vision Evaluation

Israfel Salazar[*2], Manuel Fernández Burda[*3], Shayekh Bin Islam[*4],
Arshia Soltani Moakhar[*4], Shivalika Singh[*1], Fabian Farestam[*17],
Angelika Romanou[*], Danylo Boiko[4,9], Dipika Khullar[4], Mike Zhang[10],
Dominik Krzemiński[4], Jekaterina Novikova[4], Luísa Shimabucoro[18],
Joseph Marvin Imperial[19,26], Rishabh Maheshwary[4], Sharad Duwal[4],
Alfonso Amayuelas[5], Swati Rajwal[6], Jebish Purbey[4,7], Ahmed Ruby[25],
Nicholas Popovič[11,12], Marek Suppa[22,23], Azmine Toushik Wasi[4],
Ram Mohan Rao Kadiyala[7,8], Olga Tsymboi[13, 28], Maksim Kostritsya[15,16],
Bardia Soltani Moakhar[4], Gabriel da Costa Merlin[18], Otávio Ferracioli Coletti[18],
Maral Jabbari Shiviari[4], MohammadAmin farahani fard[4], Silvia Fernandez[4],
María Grandury[21], Dmitry Abulkhanov[4], Drishti Sharma[4,7],
Andre Guarnier De Mitri[18], Leticia Bossatto Marchezi[20], Setayesh Heydari[4],
Johan Obando-Ceron[4,24], Nazar Kohut[14], Beyza Ermis[1], Desmond Elliott[♦2,27],
Enzo Ferrante[♦3], Sara Hooker[♦1], and Marzieh Fadaee[♦1]

[1]Cohere For AI, [2]Department of Computer Science, University of Copenhagen, [3]Institute of Computer Sciences, CONICET & Universidad de Buenos Aires, [4]Cohere For AI Community, [5]University of California, Santa Barbara, [6]Emory University, [7]M2ai.in, [8]Traversaal.ai, [9]Taras Shevchenko National University of Kyiv, [10]Aalborg University, [11]Karlsruhe Institute of Technology, Germany, [12]ScaDS.AI, TU Dresden, Germany, [13]T-Tech, [14]Lviv Polytechnic National University, [15]HSE University (Higher School of Economics), [16]RAFT, [17]ETH Zürich, [18]University of São Paulo, [19]National University Philippines, [20]Federal University of São Carlos, [21]SomosNLP, [22]Cisco, [23]Comenius University in Bratislava, [24]Mila, University of Montreal, [25]Uppsala University, [26]University of Bath, [27]Pioneer Center for AI, [28]Moscow Institute of Physics and Technology

Corresponding authors: Israfel Salazar <israfel.salazar@di.ku.dk>, Manuel Fernández Burda <mburda@dc.uba.ar>, Marzieh Fadaee <marzieh@cohere.com>

The evaluation of vision-language models (VLMs) has mainly relied on English-language benchmarks, leaving significant gaps in both multilingual and multicultural coverage. While multilingual benchmarks have expanded, both in size and languages, many rely on translations of English datasets, failing to capture cultural nuances. In this work, we propose KALEIDOSCOPE, as the most comprehensive exam benchmark to date for the multilingual evaluation of vision-language models. KALEIDOSCOPE is a large-scale, in-language multimodal benchmark designed to evaluate VLMs across diverse languages and visual inputs. KALEIDOSCOPE covers 18 languages and 14 different subjects, amounting to a total of 20,911 multiple-choice questions. Built through an open science collaboration with a diverse group of researchers worldwide, KALEIDOSCOPE ensures linguistic and cultural authenticity. We evaluate top-performing multilingual vision-language models and find that they perform poorly on low-resource languages and in complex multimodal scenarios. Our results highlight the need for progress on culturally inclusive multimodal evaluation frameworks.

**website:** http://cohere.com/research/kaleidoscope
**dataset:** https://hf.co/datasets/CohereForAI/kaleidoscope

---

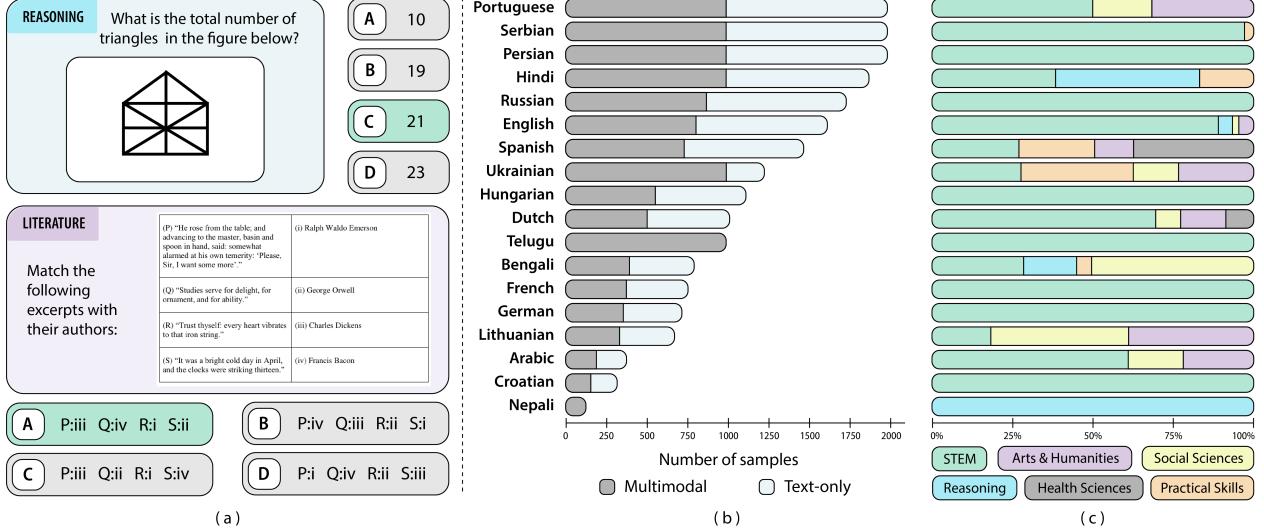[*]First authors. [♦]Principal senior advisors.

Figure 1: **Overview of the** KALEIDOSCOPE **Benchmark.** (a) Multilingual-Multimodal MCQ Samples (b) Language and Multimodal Samples Distribution. (c) Exam Category Breakdown.

# 1 Introduction

Evaluations are the backbone of measuring progress in machine learning, yet many benchmarks – especially for language models – continue to mirror an English and Western-centric worldview (Joshi et al., 2020; Fan et al., 2020; Dodge et al., 2021; Liu et al., 2021; Chung et al., 2022; Gehrmann et al., 2022; Lucy et al., 2024). This imbalance becomes even more striking at the cutting edge of AI, where generative models are rapidly expanding into multimodal territory (OpenAI et al., 2024; Google et al., 2024; Anthropic, 2024; Deitke et al., 2024; Yue et al., 2025; Qwen-Team, 2025), seeking to represent a richer world made up of different modalities such as image, text, sound. In recent years, the community has made promising strides toward broader multilingual text evaluation (Ahuja et al., 2023; Singh et al., 2024b;a; Aakanksha et al., 2024; Pozzobon et al., 2024; Romanou et al., 2024; Singh et al., 2025; Adelani et al., 2024), and multimodal benchmarks are starting to take shape (Bugliarello et al., 2022; Fu et al., 2023; Yue et al., 2024a;b; Li et al., 2024a; Xu et al., 2025). Yet reliable evaluation at the intersection of multilingual and multimodal tasks remains rare. This gap is precisely what motivates our work.

One common but imperfect solution is translating English benchmarks into other languages. While convenient, this approach often falls short of capturing cultural context and nuance. Translated datasets can easily reinforce Western-centric knowledge and assumptions (van Miltenburg et al., 2017; Frank et al., 2018; Singh et al., 2025; Longpre et al., 2025) limiting their ability to truly assess model performance across diverse settings. Moreover, automated data curation pipelines frequently amplify existing quality issues (Luccioni & Viviano, 2021; Caswell et al., 2020; Kreutzer et al., 2022), with translation artifacts such as *translationese* muddying the waters even further (Koppel & Ordan, 2011; Zhang & Toral, 2019; Bizzoni et al., 2020; Vanmassenhove et al., 2021). While translated data has its place, especially for some particularly low-resource tasks (Zhou et al., 2021; Thapliyal et al., 2022; Qiu et al., 2022; Ramos et al., 2024; Geigle et al., 2025; Dang et al., 2024; Üstün et al., 2024; Aakanksha et al., 2024), it is an imperfect substitute for genuinely diverse, in-language benchmarks.

In this work, we introduce the largest benchmark of real-world, in-language exam questions that

blend image and text modalities. Our dataset pushes beyond simple captioning tasks, challenging models to reason about visual content in various topics, the way humans are evaluated in exams worldwide. Through a large-scale open science effort across 18 languages, we construct KALEIDO-SCOPE (see Figure 1), featuring a diverse selection of knowledge domains across 14 subjects. With 55% of the total 20,911 questions requiring image understanding for accurate resolution, our work aims to establish a comprehensive, and inclusive evaluation framework for multimodal language models. We evaluate a wide range of state-of-the-art models on KALEIDOSCOPE, including Claude 3.5 Sonnet (Anthropic, 2024), GPT-4o (OpenAI et al., 2024), and Gemini-V (Google et al., 2024), as well as smaller open-weight VLMs, such as Aya-Vision model family (Cohere-For-AI-Team, 2025), Molmo (Deitke et al., 2024) Pangea (Yue et al., 2025), and Qwen2.5-VL model family (Qwen-Team, 2025). Our key contributions and findings are highlighted here:

- KALEIDOSCOPE **Benchmark**: We present the largest multilingual multimodal exam set, covering high resource (e.g., English, Spanish) to underrepresented languages (e.g., Bengali, Telugu) across diverse subjects from sociology to STEM. Most languages (10/18) include 5+ topics, with the rest focusing on multi-subtopics like mathematics or engineering. Questions emphasize vision grounded reasoning through tasks like interpreting graphs, pictures, and region-specific diagrams, supported by fine-grained metadata for model diagnostics.

- **Modality-Specific Performance Disparities:** All models perform substantially better on text-only questions, revealing a clear disparity across modalities. The gap widens in larger modelsl; for instance, GPT-4o shows a 21.6% difference between text-only and multimodal performance, while smaller models like Molmo exhibit a much narrower gap of 3.69%. (Section 4.1). Furthermore, multimodal performance varies significantly by visual data type: models are more capable of answering questions about tables (76.5%) and photographs (81.5%) compared to diagrams (62.9%).

- **Domain-Specific Performance Disparities:** We observe a significant performance gap between questions requiring knowledge of Humanities & Social Sciences and those focused on STEM subjects (Section 4.4). On average, models present accuracy of 83.7% for humanities versus 59.2% for STEM (based on the best scores across models). Models struggle more with STEM questions, suggesting that while they can often recognize visual content and retrieve related knowledge, they lack the reasoning capabilities needed to arrive at the correct answers in STEM domains.

- **Crosslingual Performance Disparities:** Model performance varies across languages, with noticeably better results in high-resource languages and weaker performance in mid- and low-resource ones (Section 4.3). Crosslingual transfer appears to play a role, as models perform better on average in languages using Latin scripts compared to those with non-Latin scripts.

## 2 The KALEIDOSCOPE Benchmark

The KALEIDOSCOPE Benchmark is a global collection of multiple-choice questions sourced from real-world exams, with the goal of evaluating multimodal and multilingual understanding in VLMs. The collected exams are in a Multiple-choice question answering (MCQA) format which provides a structured framework for evaluation by prompting models with predefined answer choices (Hendrycks et al., 2021; Lu et al., 2023; Wang et al., 2024a; Yue et al., 2024a; Romero et al., 2024; Romanou

et al., 2024), closely mimicking conventional human testing methodologies. Our work is built around three core design principles that guide the selection, curation, processing, and addition of exams:

- 🖼 **Multimodality**: Images are central to KALEIDOSCOPE, as we aim to evaluate how VLMs integrate and reason about visual information to answer questions. We prioritize multimodal questions with diverse image types, complemented by a similar proportion of text-only questions for a complete assessment and comparison.

- 🌐 **Multilinguality**: The benchmark contains questions in 18 languages, with a focus on under-represented mid- and low-resource languages (e.g., Nepali, Lithuanian) alongside high-resource languages (e.g., English, Spanish) for a thorough evaluation across a broad range of languages.

- 👥 **Diversity**: Our goal is to collect exams covering as wide a range of topics as possible ranging from Mathematics and Sociology, to Medicine and Driving Licenses, ensuring comprehensive evaluation across various domains. The final collection includes exams from 14 different domains, collected from 18 countries and with varying educational levels (from high school to professional exams), allowing detailed clustering and comprehensive evaluation.

## 2.1 Global Collaboration

Our work entailed an extensive, open science process to manually collect data by working directly with native speakers of different languages (Elliott et al., 2016; Liu et al., 2021; Thapliyal et al., 2022; Li et al., 2024c; Üstün et al., 2024; Singh et al., 2024b). This is acutely needed in the field of machine learning, where recent studies have highlighted that dataset creators remain predominantly Western-centric (Longpre et al., 2025). The manual curation of datasets is a costly process that requires careful attention to detail in every language to ensure high-quality, contextually relevant content for evaluation. In this work, we engage in a large-scale open science collection process, which brings together contributors spanning 20 nations across four continents to ensure linguistic and cultural authenticity. For related participatory research see Section 6.3.

## 2.2 Data Pipeline

**Collection:** We collected KALEIDOSCOPE with guidelines detailing information about the type of exams and questions required, formatting, specifications, and quality control measures. The data was collected through a global call for contributions and distributed across global communities, with the majority of contributors being independent researchers in the Cohere for AI (C4AI)[1] open science community. This effort resulted in a collection of 20,911 questions from 18 different countries and 18 languages, all sourced in their original languages, avoiding translations to maintain linguistic authenticity. We prioritized original, domain-expert-written questions (e.g., from teachers), ensuring real-world relevance and quality. The exams were gathered from various repositories, including official government websites, question banks, and other publicly available repositories with educational materials. Throughout the process, contributors also annotated associated licenses with each dataset to allow for documentation of data provenance (Longpre et al., 2024).

**Processing:** The annotation process consists of two stages. In the first stage, we perform automated parsing and extraction. For directly parsable text, we use PDF or web parsers, while for non-parsable

---

[1] https://cohere.com/research

| Language | Code | Subjects | Total | Visual | Text | Resources | Family |
|---|---|---|---|---|---|---|---|
| Portuguese | pt | 11 | 2000 | 1000 | 1000 | High | Italic |
| Serbian | sr | 1 | 2000 | 1000 | 1000 | High | Balto-Slavic |
| Persian | fa | 5 | 2000 | 1000 | 1000 | High | Iranian |
| Hindi | hi | 12 | 1886 | 1000 | 886 | High | Indo-Aryan |
| Russian | ru | 1 | 1744 | 872 | 872 | High | Balto-Slavic |
| English | en | 9 | 1628 | 814 | 814 | High | Germanic |
| Spanish | es | 6 | 1482 | 741 | 741 | High | Italic |
| Hungarian | hu | 1 | 1120 | 560 | 560 | High | Uralic |
| Dutch | nl | 10 | 1018 | 509 | 509 | High | Germanic |
| French | fr | 1 | 762 | 381 | 381 | High | Italic |
| German | de | 1 | 722 | 361 | 361 | High | Germanic |
| Arabic | ar | 10 | 382 | 191 | 191 | High | Semitic |
| Croatian | hr | 1 | 324 | 162 | 162 | High | Balto-Slavic |
| Ukrainian | uk | 8 | 1237 | 1000 | 237 | Mid | Balto-Slavic |
| Bengali | bn | 6 | 800 | 400 | 400 | Mid | Indo-Aryan |
| Lithuanian | lt | 6 | 680 | 340 | 340 | Mid | Balto-Slavic |
| Telugu | te | 1 | 1000 | 1000 | 0 | Low | South Dravidian |
| Nepali | ne | 1 | 126 | 126 | 0 | Low | Indo-Aryan |
| **Total** | (18) | 14 | 20,911 | 11,457 | 9,454 | – | – |

Table 1: **Statistics of the** KALEIDOSCOPE **Dataset.** Breakdown of subjects (Subjects), total questions (Total), multimodal questions (Visual), and text-only questions (Text) per language. Languages are covered by multiple sources with single-subject cases containing specialized subdomains. 🖼: Supports evaluation of both multimodal (image+text) and unimodal (text-only) capabilities. 🌐: Languages are classified by resource level (high/mid/low) following Joshi et al. (2019); Singh et al. (2024b). 👥: Enables granular analysis of model performance across modalities, languages, and subject domains.

text, we employ OCR API's, such as Mathpix[2], along with vision-language models such as GPT-4o. These tools allow us to extract both text and image elements from exam source formats, which are then converted into structured outputs in LaTeX, Markdown, and JSON formats, as required. Since automated parsing can sometimes result in misaligned images and text, the second stage involves refining the extracted text. We apply heuristic rules, as well as high-performing LLMs, such as Claude 3.5 Sonnet and GPT-4o, to restructure the output, ensuring proper alignment of questions, text, and answer choices. This stage was followed by human verification, ensuring that images are correctly linked to the corresponding questions, and checking that extracted formulas match the expected equation format.

**Quality Assessment:** Maintaining reliable and high quality data is essential, especially given the large-scale international collaboration involved in this project. To ensure data integrity, we include manual validation in three stages of the collection and annotation pipeline. First, at the end of the collection stage, two independent annotators validate each exam to ensure conformity

---

[2] https://mathpix.com/convert

with the guidelines. Part of this verification includes a strict revision to confirm compliance with the distribution license requirements. Only exams approved by both independent annotators are included in the dataset. Next, following the annotation process, a validation script checks for JSON formatting errors, duplicate questions, and malformed strings that do not conform to identified entry specifications (see Appendix A.1). Finally, at the last stage, two separate validators perform a final manual review of the collected files before merging them into KALEIDOSCOPE.

Quality control extends to the evaluation phase as well, where we analyze the most prominent failure modes. During model inference, suspicious outputs, such as ambiguous answers, no response, or consistent failures across all models, are flagged for manual review. For example, if all models fail in a specific question, we investigate further and may discover issues like missing images or incorrect labels. If an issue is identified, the entire exam that contains the problematic question is reviewed for correction or removal. This process guarantees that any errors in the benchmark questions are identified and addressed, further enhancing the reliability of the dataset.

## 2.3 Data Statistics

The final KALEIDOSCOPE benchmark contains 20,911 questions across 18 languages belonging to 8 language families. A total of 11,459 questions require an image to be answered (55%), while the remaining 9,452 (45%) are text-only. The dataset covers 14 different subjects, grouped into 6 broad domains. Figure 1 presents an overview of the dataset; detailed statistics can be found in Table 1. The majority of questions in KALEIDOSCOPE are multimodal, with the exact proportion varying across languages, ranging from 50% to 100%, with some languages always requiring images for resolution.

Each exam question contains 17 fields, including source country, language, license, educational level, category, and multimodal information. These fields are detailed in Appendix A.1. The questions are formatted in MCQA format with 4 options and a single correct answer. The subject is labeled in both English (`category_en`) and the source language (`category_source_lang`). The educational level (e.g., high school, university entrance, professional licensing) is also included to ensure diverse representation. Multimodal questions additionally specify the type of image, such as graphs, tables, or diagrams. Additionally, each entry includes metadata such as source details, licensing status, and ISO 639-1 language codes. For a fine-grained analysis, each question includes detailed metadata, with examples provided in Appendix A.2. The metadata allows us to evaluate how visual and textual elements interact in multimodal reasoning tasks, making the benchmark valuable for evaluating models across diverse scenarios.

KALEIDOSCOPE covers a wide range of languages, including low- and mid-resource languages such as Nepali, Lithuanian, Bengali, Telugu, Persian, Ukrainian, Croatian, Serbian, and Hungarian, as well as high-resource languages such as English, Spanish, Portuguese, Russian, French, German, Arabic, Hindi, and Dutch. This selection allows us to evaluate how performance is affected by the amount of resources available for a given language. The dataset spans 8 different language families, providing a broad linguistic range. The number of questions per language varies significantly, from 126 for Nepali to 2000 for Portuguese, Serbian, and Persian. The linguistic diversity present in KALEIDOSCOPE enables a robust evaluation of models across both widely spoken and underrepresented languages, making the dataset suitable for comprehensive multilingual assessment.

# 3 Experimental Setup

## 3.1 Models

We benchmark both open-weights and closed multimodal vision-language models on KALEIDO-SCOPE, focusing on lighter open-weight models and larger closed models to assess performance across a wide range of model sizes. The open-weight models[3] include Aya-Vision-8B and 32B (Cohere-For-AI-Team, 2025), Molmo-7B-D (Deitke et al., 2024), Pangea-7B (Yue et al., 2025), and all sizes of Qwen2.5-VL-Instruct (Qwen-Team, 2025) (3B, 7B, 32B, and 72B) to analyze the impact of model scale on KALEIDOSCOPE. All models have image and multilingual support; Aya-Vision supports 23 languages, Qwen2.5-VL supports 29 languages, and Pangea was trained on a dataset spanning 39 different languages, making them strong candidates for multimodal and multilingual evaluation. For the closed models, we evaluate GPT-4o (OpenAI et al., 2024) (2024/08/06), Claude 3.5 Sonnet (Anthropic, 2024) (2024/10/22), and Gemini 1.5 Pro (Google et al., 2024).

## 3.2 Evaluation Setup

We designed two distinct evaluation setups to accommodate VLMs' varying reasoning and instruction-following capabilities. For closed models, we designed zero-shot prompts using the Chain-of-Thought (CoT) method (Wei et al., 2022). The prompt instructs the model to reason through the correct answer and to write the chosen option within specific `<ANSWER> </ANSWER>` tags. This approach is natural and aligns the real-world application of MCQs, encouraging the model to generate a step-by-step reasoning before selecting the final answer. We define a common template, ensuring equal evaluation conditions for all models (see Appendix A.5). The instructions were translated to all the evaluated languages, creating a fully in-language setup, following the methodology proposed by Romanou et al. (2024). The selected choice is extracted using string matching of the tags.

For the smaller open-weight models, which typically have limited capacity for complex reasoning, CoT prompting proved less effective in our preliminary experiments (see Appendix A.7 for details). Therefore, we implemented a direct answer generation approach, instructing the models to produce a JSON output containing their choice within a predefined 'choice' field. The instruction was always in English, independent of the question language (see Appendix A.5). This setup simplifies the task, reducing errors related to multi-step reasoning or formatting inconsistencies, and ensuring a straightforward answer extraction. Further discussion about model output error analysis can be found in section 5.

## 3.3 Evaluation Metrics

Given the multiple-choice nature of the task, we use accuracy as the primary evaluation metric. We report overall accuracy across all questions, as well as accuracy on the subset of questions where the model produces valid responses. A response is considered valid if the model successfully provides an answer in the expected format and selects a valid option (i.e., one of the letters A, B, C, D). Invalid responses typically result from missing the selected choice, selecting an invalid option, or refusal to answer. To quantify these cases, we report the *Format Error Rate*, which measures the proportion

---

[3]All open-weight models are evaluated locally using 1×NVIDIA Ampere A100 GPU with 64GB of memory for models up to 8B, and 4×A100 for models on the range 32B–72B. To ensure a consistent evaluation environment, we set the temperature to 0.7, the maximum token generation to 1024, and the image size to 512×512 for all models.

| | Overall | | | Multimodal | | | Text-only | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Valid Responses | | | Valid Responses | | | Valid Responses | |
| Model | Acc. | F.E. | Acc. | Acc. | F.E. | Acc. | Acc. | F.E. | Acc. |
| Claude 3.5 Sonnet | **62.91** | 1.78 | **63.87** | **55.63** | 3.24 | **57.24** | **73.54** | 0.02 | **73.57** |
| Gemini 1.5 Pro | 62.10 | 1.62 | 62.95 | 55.01 | 1.46 | 55.71 | 72.35 | 1.81 | 73.45 |
| GPT-4o | 58.32 | 6.52 | 62.10 | 49.80 | 10.50 | 55.19 | 71.40 | 1.71 | 72.39 |
| Qwen2.5-VL-72B | 52.94 | 0.02 | 53.00 | 48.40 | 0.03 | 48.41 | 60.00 | 0.02 | 60.01 |
| Aya-Vision-32B | 39.27 | 1.05 | 39.66 | 35.74 | 1.49 | 36.28 | 44.73 | 0.51 | 45.00 |
| Qwen2.5-VL-32B | 48.21 | 0.88 | 48.64 | 44.90 | 0.28 | 45.05 | 53.77 | 1.61 | 54.60 |
| Aya-Vision-8B | 35.09 | 0.07 | 35.11 | 32.35 | 0.05 | 32.36 | 39.27 | 0.10 | 39.30 |
| Molmo-7B-D | 32.87 | 0.04 | 32.88 | 31.43 | 0.06 | 31.44 | 35.12 | 0.01 | 35.13 |
| Pangea-7B | 31.31 | 7.42 | 34.02 | 27.15 | 13.52 | 31.02 | 37.84 | 0.03 | 37.86 |
| Qwen2.5-VL-7B | 39.56 | 0.08 | 39.60 | 36.85 | 0.04 | 36.88 | 43.91 | 0.11 | 43.96 |
| Qwen2.5-VL-3B | 35.56 | 0.19 | 35.63 | 33.67 | 0.32 | 33.79 | 38.51 | 0.03 | 38.53 |

Table 2: **Performance Evaluation on** KALEIDOSCOPE. Results are reported as macro-averaged accuracy (%) across all languages (equal weight per language). **Acc.**: Accuracy over all samples; **F.E.**: Format Error rate (invalid responses); **Valid Acc.**: Accuracy excluding invalid responses. Metrics are shown for the full dataset (**Overall**), multimodal inputs (**Multimodal**), and text-only inputs (**Text-only**).

of questions for which the model fails to generate a valid answer. For grouped results, we report the macro average across languages, i.e. all languages have equal weight when computing the score.

# 4 Results

## 4.1 Overall Performance

We benchmark a wide variety of models on KALEIDOSCOPE and present the main results in Table 2. Claude 3.5 Sonnet achieves the highest overall accuracy (62.91%), followed closely by Gemini 1.5 Pro (62.10%). GPT-4o performs notably worse (58.32%), with a high format error rate (6.52% overall, with at 10.50% for the multimodal split). However, when considering only valid answers, GPT-4o's performance improves significantly (+3.78 percentage points), closing the gap with other closed models and highlighting the impact of format errors (see Section 5).

Among open-weight models, Qwen2.5-VL-72B achieves the highest accuracy (52.94%), which is expected given its larger number of parameters. In the lightweight category (≤8B parameters), Qwen2.5-VL-7B outperforms all others in both multimodal and text-only questions, with accuracy of 39.56%. Open models generally maintain low format error rates, except for Pangea, which has the highest format error rate at 13%.

Table 2 also summarizes results for both the multimodal and text-only benchmark splits. Across all models, multimodal performance is lower than text-only, with a larger drop for closed models: GPT-4o drops 21.6 accuracy points overall (10.29 points for valid answers). Open-weight models show smaller gaps, with Molmo having the narrowest gap of only 3.69 accuracy points, though multimodal performance remains lower. Among lightweight models, Qwen2.5-VL-7B leads on both splits, with a relatively small gap, followed by Aya-Vision-8B on text-only samples and Qwen2.5-
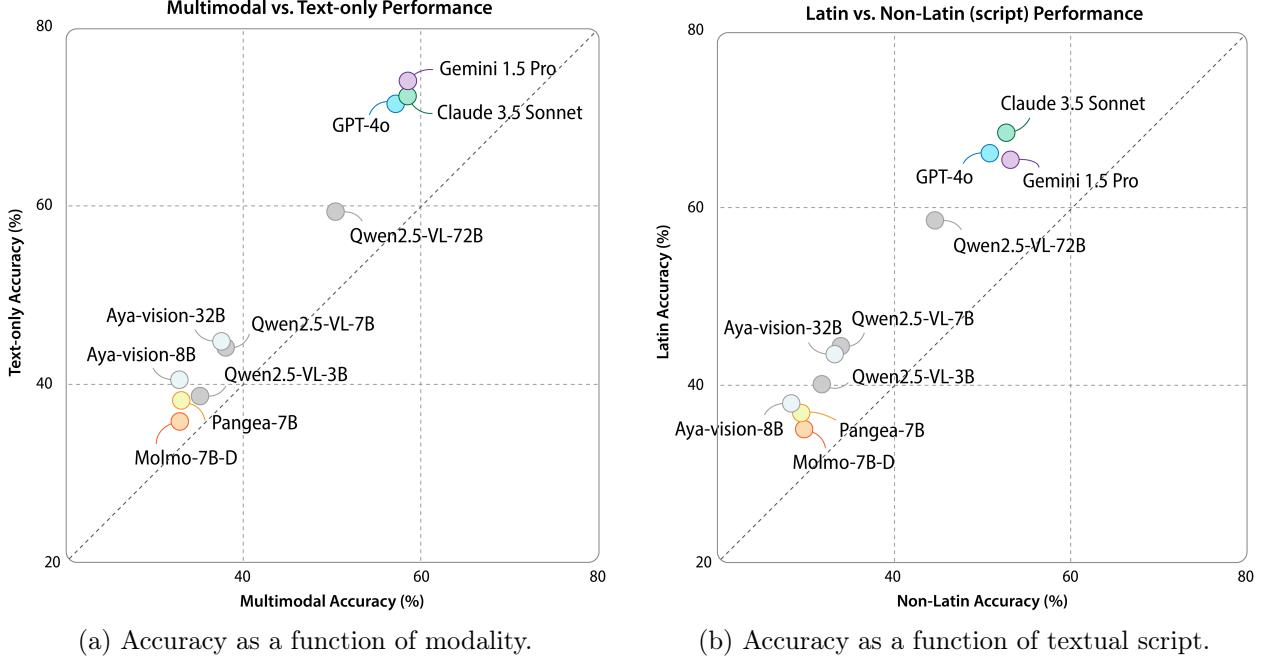
(a) Accuracy as a function of modality.

(b) Accuracy as a function of textual script.

Figure 2: **Model Performance Analysis on** KALEIDOSCOPE**.** (a) Accuracy (%) of models on multimodal and text-only questions, highlighting low performance on multimodal samples. (b) Accuracy (%) by script type, revealing biases for latin scripts. Accuracy over valid responses is used to generate both figures. Identity line is added to show parity.

VL-3B in the multimodal split. Closed models perform well on text-only questions, highlighting their strength in this modality but also the challenges of multimodal processing. In contrast, the smaller variation in performance across both splits by open-weight models suggests that they are less specialized, with lower overall performance but greater robustness in multimodal tasks. Lightweight models exhibit similar behavior across script types and modalities (Figure 2a), with Molmo showing the most balanced performance between Multimodal vs. text-only samples and Latin vs. non-Latin scripts.

## 4.2   Not All Image Types are Equal

KALEIDOSCOPE contains eight visual information types, with accuracy varying significantly by complexity (Table 3). Simpler inputs like text-rich images (Qwen2.5-VL-7B: 76.3%; GPT-4o: 86.2%) and photos score higher than technical categories like Formulas and Diagrams (Qwen2.5-VL-7B: 38.0%; GPT-4o: 62.9%). Notably, Qwen2.5-VL-72B ranks second in text-rich images, surpassing both Gemini and Claude. Larger models show specialized strengths: Gemini 1.5 Pro dominates Formulas and Figures, GPT-4o leads in text-rich images, and Claude 3.5 Sonnet achieves the highest scores in Diagrams, Graphs, and Maps. In contrast, Qwen2.5-VL-7B consistently outperforms all lightweight models across categories, demonstrating broader capability despite lower absolute scores. The results reveal a clear hierarchy: models handle simple visuals well but struggle with structured or symbolic data, a pattern consistent across architectures but more pronounced in smaller models.

9

| Model | Diagram (2,182) | Figure (6,178) | Graph (733) | Map (392) | Photo (631) | Formula (487) | Table (597) | Text (257) |
|---|---|---|---|---|---|---|---|---|
| Claude 3.5 Sonnet | **62.9** | 50.5 | **74.2** | **80.1** | 77.8 | 52.1 | 75.0 | 85.2 |
| Gemini 1.5 Pro | 59.4 | **51.3** | 67.9 | 69.4 | 75.8 | **68.3** | 76.0 | 85.2 |
| GPT-4o | 59.6 | 48.2 | 68.4 | 78.8 | **81.5** | 64.4 | **76.5** | **86.2** |
| Qwen2.5-VL-72B | 51.1 | 43.9 | 59.4 | 66.1 | 70.5 | 48.7 | 61.5 | 86.0 |
| Aya-Vision 32B | 38.6 | 33.4 | 42.0 | 50.0 | 60.2 | 32.4 | 33.1 | 68.8 |
| Qwen2.5-VL-32B | 46.7 | 41.0 | 53.1 | 58.2 | 65.0 | 47.3 | 58.0 | 82.5 |
| Aya-Vision 8B | 32.7 | 29.9 | 37.2 | 38.6 | 42.3 | 29.2 | 34.1 | 54.9 |
| Molmo-7B-D | 30.3 | 31.5 | 36.7 | 37.8 | 45.0 | 25.1 | 30.6 | 56.8 |
| Pangea-7B | 31.0 | 31.0 | 32.9 | 38.5 | 45.0 | 32.2 | 29.4 | 66.3 |
| Qwen2.5-VL-7B | 38.0 | 34.0 | 44.3 | 48.0 | 53.9 | 34.9 | 40.9 | 76.3 |
| Qwen2.5-VL-3B | 32.8 | 32.3 | 40.2 | 41.2 | 48.2 | 34.7 | 35.2 | 72.8 |

Table 3: **Model Performance Breakdown by Image Type in** KALEIDOSCOPE**.** Accuracy (%) over valid answers across image type. Bold values indicate top-performing model.

## 4.3 Resource and Script Sensitivity in Models

Performance in KALEIDOSCOPE varies widely across all 18 languages (see Figure 3). Models generally perform well in high-resource languages (e.g., English, Spanish, German) but struggle with lower- and mid-resource ones, such as Nepali and Telugu. This can be attributed to the limited training data for these languages, complex scripts, and the exclusive use of multimodal samples for these languages (see Table 1), which are inherently more challenging. Lithuanian, despite being mid-resource language, stands out as the highest-performing language, with Claude 3.5 Sonnet leading in overall accuracy. This might be due the fact that all Lithuanian questions belong to `College Graduation Exams`, and have a major subject composition of Social Sciences and Humanities in opposition to STEM subjects, which may align well with the models' capabilities. Closed models show similar performance within each language, except for German, where Claude excels. In contrast, Qwen2.5-VL-7B consistently leads all lightweight models for almost every language, and the heavier Qwen2.5-VL-72B shows the benefits of model scale.

The results show that all models are biased towards Latin script languages. As shown in Figure 2b, all models are above the parity line, exhibiting consistent higher performance for Latin scripts compared to non-Latin scripts. Full results can be found in Table 9.

## 4.4 STEM Questions Expose Model Deficiencies

KALEIDOSCOPE consists of exams covering 14 subjects and domains, with Table 4 summarizing model performance on the multimodal split across subjects. We observe that all models perform significantly better on Humanities & Social Science questions compared to other domains. The closed models achieve high accuracy in areas like Sociology (Claude: 93.4%, GPT-4o: 93.2%), Social Sciences (GPT-4o: 88.1%, Gemini: 85.7%), and Language (GPT-4o: 85.8%, Claude: 85.5%). In contrast, performance in STEM subjects, including Mathematics, Physics, and Engineering, is notably lower, with most models scoring below 50%. This suggests that while they are generally capable of recognizing visual content and retrieving surface-level knowledge, they fall short when
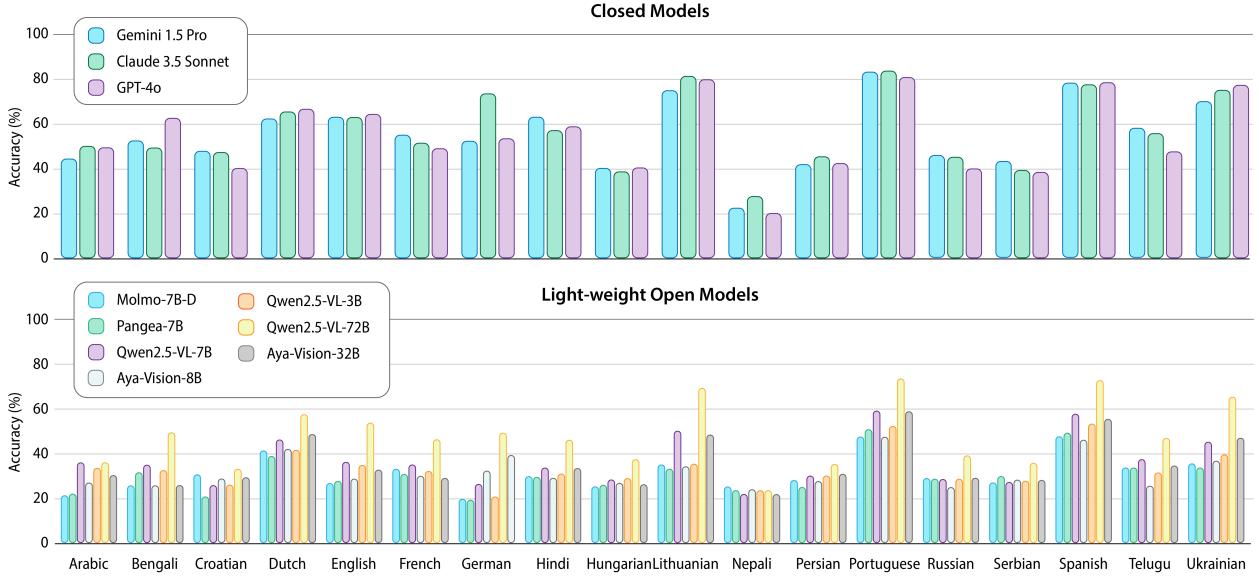
Figure 3: **Multimodal Accuracy by Language in** KALEIDOSCOPE. Reports performance (accuracy %) for closed models and open-weight models on multimodal questions.

it comes to performing the multi-step reasoning and problem-solving required in STEM subjects. Answering these questions often demands not just factual recall but also the ability to interpret complex diagrams, apply mathematical concepts, and reason through scientific principles – capabilities that current models have yet to fully master. This highlights a key gap in their ability to bridge perception and reasoning, particularly in tasks that require deeper analytical thinking.

## 5 Analysis

### 5.1 How Sensitive Are VLMs to Missing or Incorrect Images?

To evaluate the dependency of multimodal questions on images, and the impact of incorrect image associations, we conducted an experiment using the multimodal split of KALEIDOSCOPE. Following Elliott (2018); Thomason et al. (2019), we created two modified versions of the dataset: (1) a *'No Image'* split, where all images were removed, and (2) a *'Random Image'* split, where images were randomly reassigned to questions. The aim of this experiment is to assess how much the models rely on the visual information. We evaluate the performance of Qwen2.5-VL-7B on these modified splits, and the results are shown in Table 5.

We observe that the model performs above the random baseline (25%) across all three splits, indicating some ability to reason from text alone. However, there is a drop in performance ($-3.41\%$ in Total Accuracy) when questions are presented without images, suggesting that the model does rely on visual information for accurate answers. The performance drop is similar for both modifications; however, we observe a significantly larger format error when the model is tested with irrelevant images. In several of these cases, the model actually acknowledges that the image does not correspond to the question. In contrast, in experiments with no images, the format error rate is almost zero,

| | Closed Weights | | | Open Weights | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Gemini | Claude | GPT-4o | Qwen2.5-3B | Molmo-7B | Pangea-7B | Qwen2.5-7B | Aya V-8B | Aya V-32B | Qwen2.5-72B |
| *Humanities & Social Sciences* | | | | | | | | | | |
| Economics | 64.1 | 63.8 | **66.7** | 37.7 | 27.5 | 33.9 | 42.7 | 30.9 | 29.8 | 58.8 |
| Geography | 72.8 | **81.5** | 80.4 | 40.7 | 37.6 | 36.7 | 51.0 | 39.5 | 50.5 | 70.4 |
| History | 78.7 | 83.7 | **86.4** | 48.9 | 42.1 | 42.4 | 52.9 | 45.6 | 61.4 | 77.1 |
| Language | 83.5 | 85.5 | **85.8** | 72.2 | 60.1 | 66.0 | 75.7 | 56.6 | 71.2 | 85.1 |
| Social Sciences | 85.7 | 82.9 | **88.1** | 52.9 | 52.2 | 53.8 | 68.6 | 58.0 | 64.3 | 80.0 |
| Sociology | 92.3 | **93.4** | 93.2 | 64.1 | 61.0 | 57.3 | 73.1 | 57.7 | 70.5 | 87.2 |
| *STEM* | | | | | | | | | | |
| Biology | 60.3 | 62.9 | **63.9** | 37.6 | 35.3 | 33.4 | 42.6 | 35.4 | 40.7 | 53.8 |
| Chemistry | **60.4** | 59.7 | 52.9 | 33.2 | 33.5 | 34.1 | 38.5 | 28.0 | 34.8 | 50.0 |
| Engineering | 57.3 | **64.4** | 56.3 | 28.9 | 24.4 | 24.2 | 32.4 | 30.3 | 34.8 | 48.4 |
| Mathematics | **48.6** | 44.4 | 44.0 | 30.4 | 28.8 | 29.0 | 30.1 | 28.6 | 29.6 | 40.3 |
| Physics | 57.8 | **58.7** | 54.7 | 33.7 | 26.7 | 28.9 | 34.7 | 27.1 | 33.0 | 42.3 |
| *Reasoning, Health Science, and Practical Skills* | | | | | | | | | | |
| Reasoning | 52.0 | **53.3** | 51.0 | 27.4 | 27.5 | 26.6 | 29.5 | 25.1 | 27.6 | 42.3 |
| Medicine | 70.2 | 73.8 | **75.6** | 36.7 | 40.4 | 38.4 | 45.8 | 35.4 | 52.3 | 63.3 |
| Driving License | 64.4 | 64.2 | **73.1** | 39.0 | 44.9 | 39.4 | 44.9 | 41.6 | 47.1 | 54.5 |

Table 4: **Subject-wise Performance on** KALEIDOSCOPE**'s Multimodal Questions.** Valid accuracy (%) across examination subjects for multimodal samples, with bold highlighting top-performing models.

indicating that the model attempts to answer even when visual inputs are missing.[4]

## 5.2 Scaling Model Size Improves Performance

To analyze the impact of model size on KALEIDOSCOPE performance, we evaluated all four variants of Qwen2.5-VL. We selected this model family for its well-distributed size range, as well as being the best performing model in the open weight model category. We follow the same experimental setup for all model versions.

Figure 4 shows the performance of Qwen2.5-VL variants on KALEIDOSCOPE. Model size is shown in the x-axis (log-scale), while the y-axis displays accuracy for multimodal and text-only splits, and overall score. We observe a linear relationship between the logarithm of the model size and accuracy, with larger models showing significant gains. The largest open model model evaluated, Qwen2.5-VL-72B, still underperforms the closed models, however, these results highlight the effectiveness of scaling for open models, with clear and predictable improvements at each size tier.

---

[4]We observed that Qwen2.5-VL-7B tends to hallucinate when no image is present. In a simple experiment using the prompt ``Describe the following image'', the model correctly describes the input image when provided. However, when no image is passed, the model hallucinates and generates a random description.

|  | | Valid Responses | |
| Setup | Accuracy | Format Error | Accuracy |
| --- | --- | --- | --- |
| Standard Multimodal | 36.85 | 0.04 | 36.88 |
| Random Image | 32.56 | 3.12 | 33.53 |
| No Image | 33.44 | 0.03 | 33.45 |

Table 5: **Image Relevance Analysis for Qwen2.5-VL-7B on** KALEIDOSCOPE**.** Model performance across the standard multimodal, *Random Image*, and *No-Image* setups to assess the impact of visual information on question-answering accuracy.
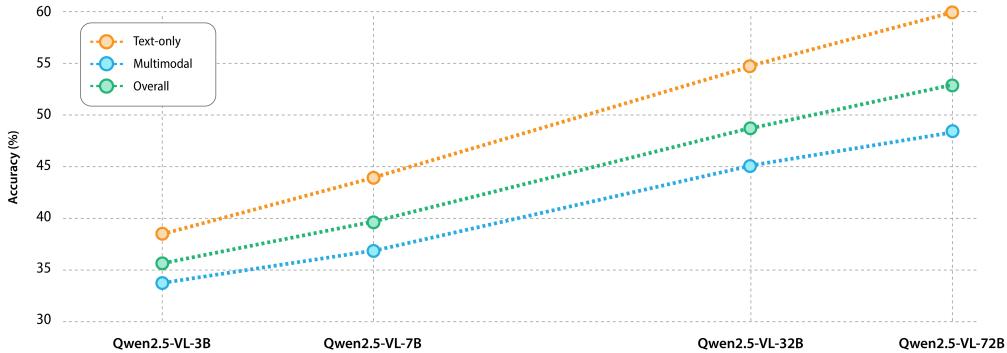


Figure 4: **Model Size Analysis for Qwen2.5-VL Models.** Performance improvement across three model sizes (3B, 7B, 32B, and 72B parameters) on KALEIDOSCOPE's multimodal tasks, demonstrating consistent gains from increased model capacity. Note that x-axis is shown in log-scale.

## 5.3 To What Extent Do Textual Augmentations Boost VLM Capabilities?

The significant performance gap between text-only and multimodal responses raises critical questions about the strengths and weaknesses of the visual processing in the tested models. In this analysis, we investigate to what extent do visual processing constraints limit multimodal capabilities, and conversely, can automatically generated textual augmentation improve model performance?

To explore this direction, we generate synthetic captions (using Gemini 1.5 Pro) and Optical Character Recognition (OCR) text (Tesseract (Smith, 2007)) for all images in KALEIDOSCOPE, aligning with the methodology of (Das et al., 2024). Unlike prior work that completely replaces images with text, we evaluate whether a VLM *augmented* with these textual inputs can boost performance.

Table 6 shows the results of augmenting visual inputs with synthetic captions and OCR text across diverse image types in KALEIDOSCOPE, measured by valid accuracy (%). Overall, the addition of a caption and OCR text improves the performance of the selected models in 5 out of 8 image types. Both models experienced a performance boost coordinately for *Graph* and *Formula*. The experiment reveals that the utility of textual augmentation depends critically on image content type. While Gemini 1.5 Pro dominates overall performance, Qwen2.5-VL-7B demonstrates selective gains when provided with captions and OCR: improvements in *Graph* (+0.9%), *Photo* (+0.2%), *Formula* (+2.4%), and *Text* (+3.5%) suggest that textual augmentation aids interpretation of content where visual elements are tightly coupled with symbolic or linguistic features (e.g., labeled axes, embedded text, or mathematical notation). Conversely, performance declines for *Diagram* (−0.1%), *Map* (−1.3%), and *Table* (−6.3%) with augmentation, implying that synthetic captions

|  | | Qwen2.5-VL-7B | | Gemini 1.5 Pro | |
|---|---|---|---|---|---|
|  | Samples | Image | +Caption | Image | +Caption |
| Diagram | 2,182 | **38.0** | 37.9 | 59.4 | **59.6** |
| Figure | 6,178 | 34.0 | **34.8** | **51.3** | 50.0 |
| Graph | 733 | 44.3 | **45.2** | 67.9 | **68.2** |
| Map | 392 | **48.0** | 46.7 | 69.4 | **70.9** |
| Photo | 631 | 53.9 | **54.1** | **75.8** | 74.3 |
| Formula | 487 | 34.9 | **37.3** | 68.3 | **68.7** |
| Table | 597 | **40.9** | 34.6 | 76.0 | **76.1** |
| Text | 257 | 76.3 | **79.8** | **85.2** | 83.7 |
| Macro Avg. | 11,457 | **36.88** | 36.83 | **55.71** | 54.81 |

Table 6: **Accuracy on augmented multimodal inputs with image captions.** Results are grouped by image type. We report **Valid Accuracy (%)**; the highest scores are highlighted in bold for each model. Macro averaged accuracy is reported over language for both methods.

may introduce noise or fail to capture structural relationships critical to these categories. Gemini's robustness across modalities ($\leq 2\%$ variation in most categories) suggests its stronger native visual understanding reduces reliance on supplementary text. The results underscore that captioning effectiveness is context-dependent: text augmentation benefits models most when (1) visual content inherently contains extractable text (e.g., *Photo* with signs, *Text* regions) or (2) symbolic patterns (e.g., formulas, graphs) require disambiguation. However, for structurally complex or text-sparse images (e.g., *Map*, *Diagram*), captioning may not compensate for deficiencies in spatial or relational reasoning. Full results, including total accuracy and format error, can be found in Table 11.

## 5.4 Format Errors

While our experimental setup ensures a majority of answers were extracted from model outputs, we observe occasional failures: models struggle to follow instructions, the outputs contain formatting errors, or models refuse to answer (particularly for health-related or ethical questions). Figure 5 shows that unanswered questions concentrate in mid- to low-resource languages, and the distribution accumulates over non-latin scripts, likely due to tokenization challenges, insufficient language-specific training data, or visual-textual alignment difficulties. Pangea-7B shows the highest refusal rates, especially for Telugu (452), Hindi (130), Persian (160), and Serbian (125). While other open models show minimal unanswered counts, indicating better format adherence. Closed models (Claude 3.5 Sonnet, GPT-4o) display distinct behavior: their refusals concentrate on non-Latin, low-resource languages, but they also show high error rates for English questions, primarily health/medical queries due to policy constraints. This underscores the trade-off between content moderation and benchmark performance.

# 6 Related Work

## 6.1 General Challenges in Multilingual Evaluation of Vision-Language Models

VLMs have demonstrated impressive performance in processing and generating text, interpreting images, and reasoning across multiple modalities. These advances have been driven by various
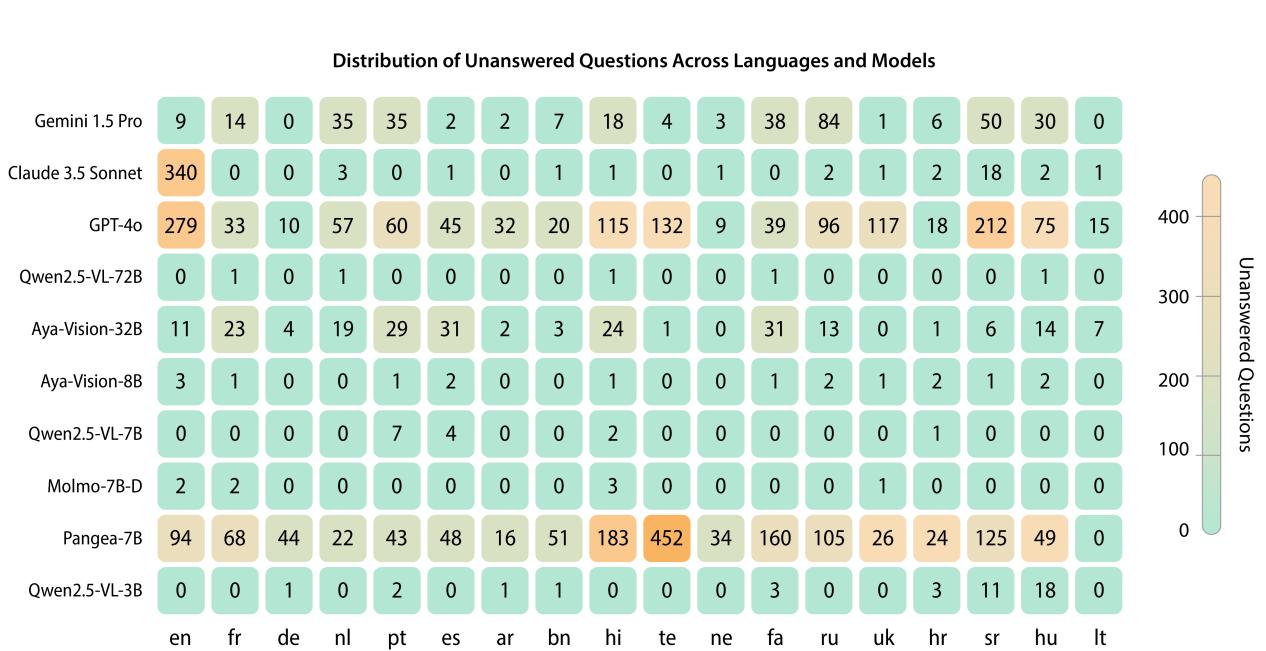
| Model | en | fr | de | nl | pt | es | ar | bn | hi | te | ne | fa | ru | uk | hr | sr | hu | lt |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Gemini 1.5 Pro | 9 | 14 | 0 | 35 | 35 | 2 | 2 | 7 | 18 | 4 | 3 | 38 | 84 | 1 | 6 | 50 | 30 | 0 |
| Claude 3.5 Sonnet | 340 | 0 | 0 | 3 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 2 | 1 | 2 | 18 | 2 | 1 |
| GPT-4o | 279 | 33 | 10 | 57 | 60 | 45 | 32 | 20 | 115 | 132 | 9 | 39 | 96 | 117 | 18 | 212 | 75 | 15 |
| Qwen2.5-VL-72B | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 |
| Aya-Vision-32B | 11 | 23 | 4 | 19 | 29 | 31 | 2 | 3 | 24 | 1 | 0 | 31 | 13 | 0 | 1 | 6 | 14 | 7 |
| Aya-Vision-8B | 3 | 1 | 0 | 0 | 1 | 2 | 0 | 0 | 1 | 0 | 0 | 1 | 2 | 1 | 2 | 1 | 2 | 0 |
| Qwen2.5-VL-7B | 0 | 0 | 0 | 0 | 7 | 4 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| Molmo-7B-D | 2 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| Pangea-7B | 94 | 68 | 44 | 22 | 43 | 48 | 16 | 51 | 183 | 452 | 34 | 160 | 105 | 26 | 24 | 125 | 49 | 0 |
| Qwen2.5-VL-3B | 0 | 0 | 1 | 0 | 2 | 0 | 1 | 1 | 0 | 0 | 0 | 3 | 0 | 0 | 3 | 11 | 18 | 0 |

Figure 5: **Distribution of the number of format errors for each model/language combination**. The languages are represented in their ISO 639 (set 1) code.

multimodal benchmarks (Li et al., 2024b; Vayani et al., 2024; Nayak et al., 2024; Schneider et al., 2025) that assess capabilities such as image captioning, object attribute recognition, and spatial relationship understanding. However, most existing benchmarks prioritize high-resource languages (e.g., English (Zang et al., 2024; Schneider et al., 2025) or Chinese (Fu et al., 2023; He et al., 2024)), resulting in significant gaps in multilingual evaluation. This focus creates disparities, particularly in evaluating performance on low-resource non-Latin languages (Hengle et al., 2024). A common strategy for lower-resource languages has been to translate existing English benchmarks using tools such as ChatGPT (Lai et al., 2023), GPT-4 (Yue et al., 2025), or Google Translate (Li et al., 2023). Such approaches, however, often fail to capture the linguistic and cultural diversity necessary for global applications (Singh et al., 2024a; Huang et al., 2025), may introduce errors or employ uncommon terms, thereby affecting the reliability of assessments, and exacerbate the problems with poverty-conscious language technology (Bird, 2022). Furthermore, the lack of comprehensive evaluation suites for non-English languages has substantially hindered multilingual generative advancements, limiting LLMs' ability to perform equitably across diverse linguistic landscapes, a challenge particularly critical for evaluating toxicity and biases in multilingual settings (Üstün et al., 2024). Addressing these limitations is essential for ensuring that VLMs perform robustly across a diverse range of real-world scenarios.

Recent efforts have attempted to address these shortcomings by incorporating culturally and linguistically diverse data. Only recently have we seen some multilingual benchmarks. For example, in the reasoning space, MMLU-ProX (Xuan et al., 2025) has created a comprehensive reasoning benchmark in 13 languages. Culturally-diverse Multilingual Visual Question Answering Benchmark (CVQA) (Romero et al., 2024) creates culturally relevant multiple-choice questions about images from 30 countries in 31 languages, using local languages which are then translated into English. In contrast, KALEIDOSCOPE builds on this work by nearly doubling the number of questions while focusing on 18 languages, thereby allowing for a more in-depth evaluation of each language. Moreover, KALEIDOSCOPE focuses on multiple-choice exams covering a wide range of topics beyond cultur-

ally relevant MCQA, including also STEM exams scenarios. Additional work in evaluating VLMs includes PangeaBench, a holistic evaluation suite encompassing 14 datasets (13 pre-existing) split between multimodal and text-only tasks, covering 47 languages (Yue et al., 2025). Similarly, Vayani et al. (2024) introduce a multimodal benchmark that includes culturally diverse images paired with text across 100 languages. Notably, this benchmark incorporates non-MCQA formats (e.g., True/False and free-form answers), which is a key distinction from the MCQA format adopted in KALEIDOSCOPE.

Other benchmarks further illustrate the diversity of evaluation approaches. For instance, MaRVL (Liu et al., 2021) assesses images in binary framework, which limits possible nuances in cultural evaluations. CULTURALVQA (Nayak et al., 2024) emphasizes cultural knowledge with approximately 44.1% of its data focusing on rituals and traditions; however, it relies on open-ended questions in English, which contrasts with the MCQA and multilingual approach of KALEIDOSCOPE. Moreover, the MaXM benchmark (Changpinyo et al., 2023) addresses bias and provides multilingual, multimodal assessment across 7 languages but does not focus on cultural aspects, an area where KALEIDOSCOPE offers added value.

## 6.2 Exam-Style Benchmarks for Vision-Language Models

Exam-style benchmarks have also advanced the evaluation of VLMs (see Table A.8). Zhang et al. (2023) present M3Exam, a novel benchmark sourced from real human exam questions that tests models in a multilingual, multimodal, and multilevel context using an MCQA framework. M3Exam includes 12,317 questions in 9 languages, requiring both multilingual proficiency and cultural knowledge; however, only about 23% of its questions necessitate image processing. This benchmark highlighted challenges faced by state-of-the-art models, such as GPT-4, particularly in handling low-resource and non-Latin script languages alongside complex multimodal queries. In contrast, KALEIDOSCOPE differentiates itself by covering a larger number of languages. Das et al. (2024) present EXAMS-V, a multi-discipline, multimodal, multilingual exam benchmark comprising 20,932 multiple-choice questions across 20 school disciplines (spanning natural sciences, social sciences, religion, fine arts, and business). EXAMS-V includes questions in 11 languages from 7 language families and incorporates four categories of multimodal features (scientific symbols, figures, graphs, and tabular data). Despite this, only 5,086 of its questions are multimodal. Additionally, the M5 benchmark (Schneider & Sitaram, 2024) evaluates VLMs on diverse vision-language tasks in a multilingual and multicultural context by covering 41 languages across eight datasets and five tasks, though it does not utilize an MCQA format. In contrast, KALEIDOSCOPE stands out from these existing exam-based benchmarks by featuring a broader diversity of languages and the largest proportion of multimodal questions in an MCQA format, with 55% of its questions requiring image understanding. This makes KALEIDOSCOPE, a comprehensive and challenging multimodal and multilingual testing ground for evaluating vision-language models in multilingual real-scenarios.

## 6.3 Participatory Open Science Projects

Participatory research empowers diverse communities to actively contribute to research processes, capturing linguistic subtleties and cultural nuances directly from native speakers. Prior participatory NLP research has primarily targeted region-specific tasks such as translation, character recognition, and audio transcription. We highlight notable initiatives here which served as our motivation and backbone framework for building KALEIDOSCOPE.

In Africa, the **Masakhane**[5] community exemplifies impactful participatory NLP by focusing on grassroots-led data collection, annotation, and model creation for African languages. Nekoto et al. (2020) demonstrated that communities in low-resource environments significantly contribute to NLP, even without formal training. Subsequent efforts by Adelani et al. (2023) have further advanced dataset curation and model development for underrepresented African languages using similar participatory frameworks. Similarly, the **MaRVL** dataset (Multicultural Reasoning over Vision and Language; Liu et al., 2021) employed native speakers from diverse linguistic backgrounds (*Indonesian, Swahili, Tamil, Turkish*, and *Mandarin Chinese*) to contribute culturally representative images, subsequently annotated by professional linguists. Despite its cultural richness, MaRVL's modest scale (under 8,000 data points) limits broader applicability beyond evaluation.

In Latin America, participatory research has also emerged and is continuously growing through the help of communities. Recent works include Hernandez Mena & Meza Ruiz (2022), which developed eight open-access linguistic resources via structured social service programs, engaging student volunteers in transcription and segmentation tasks. Concurrently, Cañete et al. (2020) and Guevara-Rukoz et al. (2020) spearheaded crowd-sourced corpora addressing dialectal diversity and resource scarcity specific to Latin American Spanish.

In Southeast Asia, **Project SEALD**[6], a collaboration between AI Singapore and Google Research, facilitated multilingual dataset collection to support regional Large Language Models (LLMs). Outputs from SEALD underpin open-source multilingual models such as *SEA-LION*[7], *Wangchan-Lion* (Phatthiyaphaibun et al., 2024), and *Sahabat-AI*[8]. Related initiatives include **NusaCrowd** for aggregating and standardizing Indonesian NLP datasets (Cahyawijaya et al., 2023) and the **SEACrowd** and **SEA-VL** projects aimed at comprehensive evaluation and benchmarking of LLMs across Southeast Asian languages (Cahyawijaya et al., 2025; Lovenia et al., 2024).

On a global scale, the **CVQA dataset** (Romero et al., 2024) was created using a participatory approach, involving native speakers and cultural experts from over 30 countries. Annotators were selected for their fluency in local languages and cultural familiarity. Many contributors were also recognized as co-authors based on their level of involvement, reinforcing a collaborative, community-driven effort. The **Aya Initiative** employed participatory methods, engaging over 3,000 contributors to curate instruction datasets across 114 languages, resulting in one of the largest multilingual datasets for language model training (Singh et al., 2024b; Üstün et al., 2024). Similarly, the IN-CLUDE benchmark (Romanou et al., 2024) leveraged participatory approaches closely aligned with our methodology. The **BigScience ROOTS corpus**, developed collaboratively for the BLOOM model, exemplifies large-scale participatory data collection. Approximately 62% of ROOTS data was crowd-sourced via global hackathons and open submissions, involving over 1,000 researchers from 60 countries and more than 250 institutions, resulting in 1.6 terabytes of multilingual data (Laurençon et al., 2022). Additionally, Uzuner et al. (2010) underscored the viability of community-driven annotation for complex, domain-specific NLP tasks like clinical text annotation, highlighting broader applicability of participatory frameworks beyond general NLP domains.

Participatory methods have also successfully extended into reinforcement learning from human

---

[5] https://www.masakhane.io/

[6] **S**outh**e**ast **A**sian **L**anguages in One Network **D**ata; https://aisingapore.org/aiproducts/southeast-asian-languages-in-one-network-data-seald/

[7] https://sea-lion.ai

[8] https://sahabat-ai.com

feedback (RLHF). For instance, the **OpenAssistant** project, led by LAION, utilized global crowd-sourcing to construct a multilingual corpus comprising over 161,000 messages annotated by 13,500 volunteers. This dataset facilitated robust training of dialogue-aligned language models through extensive human feedback annotations (Köpf et al., 2023).

## 7 Conclusion

As generative models become increasingly multimodal and multilingual, the need for robust and culturally grounded evaluation benchmarks has never been more urgent. In this work, we take a step toward closing this gap by introducing the largest benchmark of real-world, in-language multimodal exam questions. By grounding evaluation in authentic exam settings from around the world, our benchmark challenges models to reason about images in ways that mirror human assessment, capturing both linguistic and cultural complexity.

Our findings highlight the limitations of current models in handling this intersection of skills: multilingual understanding, visual reasoning, and culturally aware problem-solving. We hope this benchmark serves not only as a valuable tool for measuring progress but also as a call to action for developing models that are truly capable of operating across languages, cultures, and modalities. Continued investment in representative, high-quality evaluation datasets will be essential to ensure that future AI systems are equitable and globally relevant.

## Limitations

While our benchmark represents an important step toward more representative multilingual multimodal evaluations, several limitations still remain. First, the dataset is inherently imbalanced across languages. Coverage varies depending on the availability and accessibility of exam sources, with some languages significantly underrepresented. Second, difficulty levels are not uniformly controlled. Since questions are drawn directly from real-world exams across diverse educational systems, variations in exam design, curricular focus, and intended grade levels introduce potential inconsistency in task complexity across languages and modalities. Further the chosen MCQA question format, inherent to many exams, has issues, see Appendix B. For instance: **Exploitation of biases:** Models may guess correct answers by exploiting statistical patterns or poorly designed distractors, inflating performance metrics without demonstrating genuine understanding. **Limited real-world applicability:** Unlike open-ended queries typical in real-world applications, MCQA provides predefined options, which may not reflect natural user interactions. **Choice-order sensitivity:** Performance can vary based on the order of answer choices, introducing inconsistencies unrelated to model capability. Finally, while the dataset expands coverage beyond English, the overall language diversity remains limited. Many languages, especially those spoken in low-resource regions, are still missing due to the scarcity of suitable exam material and annotators.

## Acknowledgments

# References

Aakanksha, Arash Ahmadian, Beyza Ermis, Seraphina Goldfarb-Tarrant, Julia Kreutzer, Marzieh Fadaee, and Sara Hooker. The multilingual alignment prism: Aligning global and local preferences to reduce harm. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 12027–12049, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653 /v1/2024.emnlp-main.671. URL https://aclanthology.org/2024.emnlp-main.671/.

David Ifeoluwa Adelani, Marek Masiak, Israel Abebe Azime, Jesujoba Alabi, Atnafu Lambebo Tonja, Christine Mwase, Odunayo Ogundepo, Bonaventure F. P. Dossou, Akintunde Oladipo, Doreen Nixdorf, Chris Chinenye Emezue, Sana Al-azzawi, Blessing Sibanda, Davis David, Lolwethu Ndolela, Jonathan Mukiibi, Tunde Ajayi, Tatiana Moteu, Brian Odhiambo, Abraham Owodunni, Nnaemeka Obiefuna, Muhidin Mohamed, Shamsuddeen Hassan Muhammad, Teshome Mulugeta Ababu, Saheed Abdullahi Salahudeen, Mesay Gemeda Yigezu, Tajuddeen Gwadabe, Idris Abdulmumin, Mahlet Taye, Oluwabusayo Awoyomi, Iyanuoluwa Shode, Tolulope Adelani, Habiba Abdulganiyu, Abdul-Hakeem Omotayo, Adetola Adeeko, Abeeb Afolabi, Anuoluwapo Aremu, Olanrewaju Samuel, Clemencia Siro, Wangari Kimotho, Onyekachi Ogbu, Chinedu Mbonu, Chiamaka Chukwuneke, Samuel Fanijo, Jessica Ojo, Oyinkansola Awosan, Tadesse Kebede, Toadoum Sari Sakayo, Pamela Nyatsine, Freedmore Sidume, Oreen Yousuf, Mardiyyah Oduwole, Kanda Tshinu, Ussen Kimanuka, Thina Diko, Siyanda Nxakama, Sinodos Nigusse, Abdulmejid Johar, Shafie Mohamed, Fuad Mire Hassan, Moges Ahmed Mehamed, Evrard Ngabire, Jules Jules, Ivan Ssenkungu, and Pontus Stenetorp. MasakhaNEWS: News topic classification for African languages. In Jong C. Park, Yuki Arase, Baotian Hu, Wei Lu, Derry Wijaya, Ayu Purwarianti, and Adila Alfa Krisnadhi (eds.), *Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 144–159, Nusa Dua, Bali, November 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.ijcnlp-main.10. URL https://aclanthology.org/2023.ijcnlp-main.10/.

David Ifeoluwa Adelani, Jessica Ojo, Israel Abebe Azime, Zhuang Yun Jian, Jesujoba Oluwadara Alabi, Xuanli He, Millicent Ochieng, Sara Hooker, Andiswa Bukula, En-Shiun Annie Lee, Chiamaka Chukwuneke, Happy Buzaaba, Blessing K. Sibanda, Godson Kalipe, Jonathan Mukiibi, Salomon Kabongo KABENAMUALU, Foutse Yuehgoh, Mmasibidi Setaka, Lolwethu Ndolela, Nkiruka Bridget Odu, Rooweither Mabuya, Shamsuddeen Hassan Muhammad, Salomey Osei, Sokhar Samb, Tadesse Kebede Guge, and Pontus Stenetorp. Irokobench: A new benchmark for african languages in the age of large language models. *ArXiv*, abs/2406.03368, 2024. URL https://api.semanticscholar.org/CorpusID:270258352.

Kabir Ahuja, Harshita Diddee, Rishav Hada, Millicent Ochieng, Krithika Ramesh, Prachi Jain, Akshay Nambi, Tanuja Ganu, Sameer Segal, Maxamed Axmed, Kalika Bali, and Sunayana Sitaram. Mega: Multilingual evaluation of generative ai, 2023. URL https://arxiv.org/abs/2303.12528.

Anthropic. The Claude 3 model family: Opus, sonnet, haiku, 2024. URL https://api.semantic scholar.org/CorpusID:268232499.

Steven Bird. Local languages, third spaces, and other high-resource scenarios. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 7817–7829, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long. 539. URL https://aclanthology.org/2022.acl-long.539/.

Yuri Bizzoni, Tom S Juzek, Cristina España-Bonet, Koel Dutta Chowdhury, Josef van Genabith, and Elke Teich. How human is machine translationese? comparing human and machine translations of text and speech. In Marcello Federico, Alex Waibel, Kevin Knight, Satoshi Nakamura, Hermann Ney, Jan Niehues, Sebastian Stüker, Dekai Wu, Joseph Mariani, and Francois Yvon (eds.), *Proceedings of the 17th International Conference on Spoken Language Translation*, pp. 280–290, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.iwslt-1.34. URL https://aclanthology.org/2020.iwslt-1.34/.

Emanuele Bugliarello, Fangyu Liu, Jonas Pfeiffer, Siva Reddy, Desmond Elliott, Edoardo Maria Ponti, and Ivan Vulić. IGLUE: A benchmark for transfer learning across modalities, tasks, and languages. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 2370–2392. PMLR, 17–23 Jul 2022. URL https://proceedings.mlr.press/v162/bugliarello22a.html.

Samuel Cahyawijaya, Holy Lovenia, Alham Fikri Aji, Genta Winata, Bryan Wilie, Fajri Koto, Rahmad Mahendra, Christian Wibisono, Ade Romadhony, Karissa Vincentio, Jennifer Santoso, David Moeljadi, Cahya Wirawan, Frederikus Hudi, Muhammad Satrio Wicaksono, Ivan Parmonangan, Ika Alfina, Ilham Firdausi Putra, Samsul Rahmadani, Yulianti Oenang, Ali Septiandri, James Jaya, Kaustubh Dhole, Arie Suryani, Rifki Afina Putri, Dan Su, Keith Stevens, Made Nindyatama Nityasya, Muhammad Adilazuarda, Ryan Hadiwijaya, Ryandito Diandaru, Tiezheng Yu, Vito Ghifari, Wenliang Dai, Yan Xu, Dyah Damapuspita, Haryo Wibowo, Cuk Tho, Ichwanul Karo Karo, Tirana Fatyanosa, Ziwei Ji, Graham Neubig, Timothy Baldwin, Sebastian Ruder, Pascale Fung, Herry Sujaini, Sakriani Sakti, and Ayu Purwarianti. NusaCrowd: Open source initiative for Indonesian NLP resources. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 13745–13818, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023 .findings-acl.868. URL https://aclanthology.org/2023.findings-acl.868/.

Samuel Cahyawijaya, Holy Lovenia, Joel Ruben Antony Moniz, Tack Hwa Wong, Mohammad Rifqi Farhansyah, Thant Thiri Maung, Frederikus Hudi, David Anugraha, Muhammad Ravi Shulthan Habibi, Muhammad Reza Qorib, et al. Crowdsource, crawl, or generate? creating sea-vl, a multicultural vision-language dataset for southeast asia. *arXiv preprint arXiv:2503.07920*, 2025. URL https://arxiv.org/abs/2503.07920.

Isaac Caswell, Theresa Breiner, Daan van Esch, and Ankur Bapna. Language ID in the wild: Unexpected challenges on the path to a thousand-language web text corpus. In Donia Scott, Nuria Bel, and Chengqing Zong (eds.), *Proceedings of the 28th International Conference on Computational Linguistics*, pp. 6588–6608, Barcelona, Spain (Online), December 2020. International Committee on Computational Linguistics. doi: 10.18653/v1/2020.coling-main.579. URL https://aclanthology.org/2020.coling-main.579/.

José Cañete, Gabriel Chaperon, Rodrigo Fuentes, Jou-Hui Ho, Hojin Kang, and Jorge Pérez. Spanish pre-trained bert model and evaluation data. In *PML4DC at ICLR 2020*, 2020.

Beer Changpinyo, Linting Xue, Michal Yarom, Ashish Thapliyal, Idan Szpektor, Julien Amelot, Xi Chen, and Radu Soricut. Maxm: Towards multilingual visual question answering. In *Findings of ACL: EMNLP*, 2023. URL https://arxiv.org/abs/2209.05401.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language models, 2022. URL https://arxiv.org/abs/2210.11416.

Cohere-For-AI-Team. Aya vision: Expanding the worlds ai can see, March 2025. URL https://cohere.com/blog/aya-vision.

John Dang, Arash Ahmadian, Kelly Marchisio, Julia Kreutzer, Ahmet Üstün, and Sara Hooker. RLHF can speak many languages: Unlocking multilingual preference optimization for LLMs. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 13134–13156, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main. 729. URL https://aclanthology.org/2024.emnlp-main.729/.

Rocktim Jyoti Das, Simeon Emilov Hristov, Haonan Li, Dimitar Iliyanov Dimitrov, Ivan Koychev, and Preslav Nakov. EXAMS-V: A multi-discipline multilingual multimodal exam benchmark for evaluating vision language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 7768–7791, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.420. URL https://aclanthology.org/2024.acl-long.420/.

Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, Jiasen Lu, Taira Anderson, Erin Bransom, Kiana Ehsani, Huong Ngo, YenSung Chen, Ajay Patel, Mark Yatskar, Chris Callison-Burch, Andrew Head, Rose Hendrix, Favyen Bastani, Eli VanderBilt, Nathan Lambert, Yvonne Chou, Arnavi Chheda, Jenna Sparks, Sam Skjonsberg, Michael Schmitz, Aaron Sarnat, Byron Bischoff, Pete Walsh, Chris Newell, Piper Wolters, Tanmay Gupta, Kuo-Hao Zeng, Jon Borchardt, Dirk Groeneveld, Crystal Nam, Sophie Lebrecht, Caitlin Wittlif, Carissa Schoenick, Oscar Michel, Ranjay Krishna, Luca Weihs, Noah A. Smith, Hannaneh Hajishirzi, Ross Girshick, Ali Farhadi, and Aniruddha Kembhavi. Molmo and pixmo: Open weights and open data for state-of-the-art vision-language models. *arXiv preprint arXiv:2409.17146*, 2024. URL https://arxiv.org/abs/2409.17146.

Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld, Margaret Mitchell, and Matt Gardner. Documenting large webtext corpora: A case study on the colossal clean crawled corpus. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 1286–1305, Online and Punta Cana, Dominican Republic, November

2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.98. URL https://aclanthology.org/2021.emnlp-main.98/.

Mengnan Du, Fengxiang He, Na Zou, Dacheng Tao, and Xia Hu. Shortcut learning of large language models in natural language understanding. *Commun. ACM*, 67(1):110–120, December 2023. ISSN 0001-0782. doi: 10.1145/3596490. URL https://doi.org/10.1145/3596490.

Desmond Elliott. Adversarial evaluation of multimodal machine translation. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun'ichi Tsujii (eds.), *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 2974–2978, Brussels, Belgium, October-November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1329. URL https://aclanthology.org/D18-1329/.

Desmond Elliott, Stella Frank, Khalil Sima'an, and Lucia Specia. Multi30K: Multilingual English-German image descriptions. In Anya Belz, Erkut Erdem, Krystian Mikolajczyk, and Katerina Pastra (eds.), *Proceedings of the 5th Workshop on Vision and Language*, pp. 70–74, Berlin, Germany, August 2016. Association for Computational Linguistics. doi: 10.18653/v1/W16-3210. URL https://aclanthology.org/W16-3210/.

Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, Naman Goyal, Tom Birch, Vitaliy Liptchinsky, Sergey Edunov, Edouard Grave, Michael Auli, and Armand Joulin. Beyond english-centric multilingual machine translation, 2020. URL https://arxiv.org/abs/2010.11125.

Stella Frank, Desmond Elliott, and Lucia Specia. Assessing multilingual multimodal image description: Studies of native speaker preferences and translator choices. *Natural Language Engineering*, 24(3):393–413, 2018. doi: 10.1017/S1351324918000074.

Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023.

Sebastian Gehrmann, Abhik Bhattacharjee, Abinaya Mahendiran, Alex Wang, Alexandros Papangelis, Aman Madaan, Angelina Mcmillan-major, Anna Shvets, Ashish Upadhyay, Bernd Bohnet, Bingsheng Yao, Bryan Wilie, Chandra Bhagavatula, Chaobin You, Craig Thomson, Cristina Garbacea, Dakuo Wang, Daniel Deutsch, Deyi Xiong, Di Jin, Dimitra Gkatzia, Dragomir Radev, Elizabeth Clark, Esin Durmus, Faisal Ladhak, Filip Ginter, Genta Indra Winata, Hendrik Strobelt, Hiroaki Hayashi, Jekaterina Novikova, Jenna Kanerva, Jenny Chim, Jiawei Zhou, Jordan Clive, Joshua Maynez, João Sedoc, Juraj Juraska, Kaustubh Dhole, Khyathi Raghavi Chandu, Laura Perez Beltrachini, Leonardo F . R. Ribeiro, Lewis Tunstall, Li Zhang, Mahim Pushkarna, Mathias Creutz, Michael White, Mihir Sanjay Kale, Moussa Kamal Eddine, Nico Daheim, Nishant Subramani, Ondrej Dusek, Paul Pu Liang, Pawan Sasanka Ammanamanchi, Qi Zhu, Ratish Puduppully, Reno Kriz, Rifat Shahriyar, Ronald Cardenas, Saad Mahamood, Salomey Osei, Samuel Cahyawijaya, Sanja Štajner, Sebastien Montella, Shailza Jolly, Simon Mille, Tahmid Hasan, Tianhao Shen, Tosin Adewumi, Vikas Raunak, Vipul Raheja, Vitaly Nikolaev, Vivian Tsai, Yacine Jernite, Ying Xu, Yisi Sang, Yixin Liu, and Yufang Hou. GEMv2: Multilingual NLG benchmarking in a single line of code. In Wanxiang Che and Ekaterina Shutova (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 266–281, Abu Dhabi, UAE, December 2022.

Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-demos.27. URL https://aclanthology.org/2022.emnlp-demos.27/.

Gregor Geigle, Florian Schneider, Carolin Holtermann, Chris Biemann, Radu Timofte, Anne Lauscher, and Goran Glavaš. Centurio: On drivers of multilingual ability of large vision-language model, 2025. URL https://arxiv.org/abs/2501.05122.

Gemini Team Google, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context, 2024. URL https://arxiv.org/abs/2403.05530.

Adriana Guevara-Rukoz, Isin Demirsahin, Fei He, Shan-Hui Cathy Chu, Supheakmungkol Sarin, Knot Pipatsrisawat, Alexander Gutkin, Alena Butryna, and Oddur Kjartansson. Crowdsourcing Latin American Spanish for low-resource text-to-speech. In Nicoletta Calzolari, Frédéric Béchet, Philippe Blache, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, Hélène Mazo, Asuncion Moreno, Jan Odijk, and Stelios Piperidis (eds.), *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pp. 6504–6513, Marseille, France, May 2020. European Language Resources Association. ISBN 979-10-95546-34-4. URL https://aclanthology.org/2020.lrec-1.801/.

Zheqi He, Xinya Wu, Pengfei Zhou, Richeng Xuan, Guang Liu, Xi Yang, Qiannan Zhu, and Hua Huang. Cmmu: A benchmark for chinese multi-modal multi-type question understanding and reasoning. *arXiv preprint arXiv:2401.14011*, 2024.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=d7KBjmI3GmQ.

Amey Hengle, Prasoon Bajpai, Soham Dan, and Tanmoy Chakraborty. Multilingual needle in a haystack: Investigating long-context behavior of multilingual large language models, 2024. URL https://arxiv.org/abs/2408.10151.

Carlos Daniel Hernandez Mena and Ivan Vladimir Meza Ruiz. Creating Mexican Spanish language resources through the social service program. In Chris Callison-Burch, Christopher Cieri, James Fiumara, and Mark Liberman (eds.), *Proceedings of the 2nd Workshop on Novel Incentives in Data Collection from People: models, implementations, challenges and results within LREC 2022*, pp. 20–24, Marseille, France, June 2022. European Language Resources Association. URL https://aclanthology.org/2022.nidcp-1.4/.

Kaiyu Huang, Fengran Mo, Xinyu Zhang, Hongliang Li, You Li, Yuanchi Zhang, Weijian Yi, Yulong Mao, Jinchen Liu, Yuzhuang Xu, Jinan Xu, Jian-Yun Nie, and Yang Liu. A survey on large language models with multilingualism: Recent advances and new frontiers, 2025. URL https://arxiv.org/abs/2405.10936.

Pratik Joshi, Christain Barnes, Sebastin Santy, Simran Khanuja, Sanket Shah, Anirudh Srinivasan, Satwik Bhattamishra, Sunayana Sitaram, Monojit Choudhury, and Kalika Bali. Unsung challenges of building and deploying language technologies for low resource language communities. In Dipti Misra Sharma and Pushpak Bhattacharya (eds.), *Proceedings of the 16th International Conference on Natural Language Processing*, pp. 211–219, International Institute of Information Technology, Hyderabad, India, December 2019. NLP Association of India. URL https://aclanthology.org/2019.icon-1.25/.

Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. The state and fate of linguistic diversity and inclusion in the NLP world. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 6282–6293, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.560. URL https://aclanthology.org/2020.acl-main.560/.

Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi Rui Tam, Keith Stevens, Abdullah Barhoum, Duc Nguyen, Oliver Stanley, Richárd Nagyfi, Shahul ES, Sameer Suri, David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and Alexander Mattick. Openassistant conversations - democratizing large language model alignment. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 47669–47681. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/949f0f8f32267d297c2d4e3ee10a2e7e-Paper-Datasets_and_Benchmarks.pdf.

Moshe Koppel and Noam Ordan. Translationese and its dialects. In Dekang Lin, Yuji Matsumoto, and Rada Mihalcea (eds.), *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pp. 1318–1326, Portland, Oregon, USA, June 2011. Association for Computational Linguistics. URL https://aclanthology.org/P11-1132/.

Julia Kreutzer, Isaac Caswell, Lisa Wang, Ahsan Wahab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, Monang Setyawan, Supheakmungkol Sarin, Sokhar Samb, Benoît Sagot, Clara Rivera, Annette Rios, Isabel Papadimitriou, Salomey Osei, Pedro Ortiz Suarez, Iroro Orife, Kelechi Ogueji, Andre Niyongabo Rubungo, Toan Q. Nguyen, Mathias Müller, André Müller, Shamsuddeen Hassan Muhammad, Nanda Muhammad, Ayanda Mnyakeni, Jamshidbek Mirzakhalov, Tapiwanashe Matangira, Colin Leong, Nze Lawson, Sneha Kudugunta, Yacine Jernite, Mathias Jenny, Orhan Firat, Bonaventure F. P. Dossou, Sakhile Dlamini, Nisansa de Silva, Sakine Çabuk Ballı, Stella Biderman, Alessia Battisti, Ahmed Baruwa, Ankur Bapna, Pallavi Baljekar, Israel Abebe Azime, Ayodele Awokoya, Duygu Ataman, Orevaoghene Ahia, Oghenefego Ahia, Sweta Agrawal, and Mofetoluwa Adeyemi. Quality at a glance: An audit of web-crawled multilingual datasets. *Transactions of the Association for Computational Linguistics*, 10:50–72, 2022. doi: 10.1162/tacl_a_00447. URL https://aclanthology.org/2022.tacl-1.4/.

Viet Lai, Chien Nguyen, Nghia Ngo, Thuat Nguyen, Franck Dernoncourt, Ryan Rossi, and Thien Nguyen. Okapi: Instruction-tuned large language models in multiple languages with reinforcement learning from human feedback. In Yansong Feng and Els Lefever (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 318–327, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-demo.28. URL https://aclanthology.org/2023.emnlp-demo.28/.

Hugo Laurençon, Lucile Saulnier, Thomas Wang, Christopher Akiki, Albert Villanova del Moral, Teven Le Scao, Leandro Von Werra, Chenghao Mou, Eduardo González Ponferrada, Huu Nguyen, Jörg Frohberg, Mario Šaško, Quentin Lhoest, Angelina McMillan-Major, Gerard Dupont, Stella Biderman, Anna Rogers, Loubna Ben allal, Francesco De Toni, Giada Pistilli, Olivier Nguyen, Somaieh Nikpoor, Maraim Masoud, Pierre Colombo, Javier de la Rosa, Paulo Villegas, Tristan Thrush, Shayne Longpre, Sebastian Nagel, Leon Weber, Manuel Muñoz, Jian Zhu, Daniel Van Strien, Zaid Alyafeai, Khalid Almubarak, Minh Chien Vu, Itziar Gonzalez-Dios, Aitor Soroa, Kyle Lo, Manan Dey, Pedro Ortiz Suarez, Aaron Gokaslan, Shamik Bose, David Adelani, Long

Phan, Hieu Tran, Ian Yu, Suhas Pai, Jenny Chim, Violette Lepercq, Suzana Ilic, Margaret Mitchell, Sasha Alexandra Luccioni, and Yacine Jernite. The bigscience roots corpus: A 1.6tb composite multilingual dataset. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 31809–31826. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/ce9e92e3de2372a4b93353eb7f3dc0bd-Paper-Datasets_and_Benchmarks.pdf.

Bohao Li, Yuying Ge, Yixiao Ge, Guangzhi Wang, Rui Wang, Ruimao Zhang, and Ying Shan. Seed-bench: Benchmarking multimodal large language models. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 13299–13308, 2024a. doi: 10.1109/CVPR52733.2024.01263.

Haonan Li, Fajri Koto, Minghao Wu, Alham Fikri Aji, and Timothy Baldwin. Bactrian-x: Multilingual replicable instruction-following models with low-rank adaptation. *arXiv preprint arXiv:2305.15011*, 2023.

Jian Li, Weiheng Lu, Hao Fei, Meng Luo, Ming Dai, Min Xia, Yizhang Jin, Zhenye Gan, Ding Qi, Chaoyou Fu, et al. A survey on benchmarks of multimodal large language models. *arXiv preprint arXiv:2408.08632*, 2024b.

Wenyan Li, Crystina Zhang, Jiaang Li, Qiwei Peng, Raphael Tang, Li Zhou, Weijia Zhang, Guimin Hu, Yifei Yuan, Anders Søgaard, Daniel Hershcovich, and Desmond Elliott. FoodieQA: A multimodal dataset for fine-grained understanding of Chinese food culture. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 19077–19095, Miami, Florida, USA, November 2024c. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.1063. URL https://aclanthology.org/2024.emnlp-main.1063/.

Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pp. 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics. URL https://aclanthology.org/W04-1013/.

Fangyu Liu, Emanuele Bugliarello, Edoardo Maria Ponti, Siva Reddy, Nigel Collier, and Desmond Elliott. Visually grounded reasoning across languages and cultures. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 10467–10485, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.818. URL https://aclanthology.org/2021.emnlp-main.818/.

Shayne Longpre, Robert Mahari, Anthony Chen, Naana Obeng-Marnu, Damien Sileo, William Brannon, Niklas Muennighoff, Nathan Khazam, Jad Kabbara, Kartik Perisetla, Xinyi Wu, Enrico Shippole, Kurt Bollacker, Tongshuang Wu, Luis Villa, Alex Pentland, and Sara Hooker. A large-scale audit of dataset licensing and attribution in ai. *Nature Machine Intelligence*, 6:975–987, 08 2024. doi: 10.1038/s42256-024-00878-8.

Shayne Longpre, Nikhil Singh, Manuel Cherep, Kushagra Tiwary, Joanna Materzynska, William Brannon, Robert Mahari, Naana Obeng-Marnu, Manan Dey, Mohammed Hamdy, Nayan Saxena, Ahmad Mustafa Anis, Emad A. Alghamdi, Vu Minh Chien, Da Yin, Kun Qian, Yizhi Li, Minnie Liang, An Dinh, Shrestha Mohanty, Deividas Mataciunas, Tobin South, Jianguo Zhang, Ariel N. Lee, Campbell S. Lund, Christopher Klamm, Damien Sileo, Diganta Misra, Enrico Shippole, Kevin Klyman, Lester JV Miranda, Niklas Muennighoff, Seonghyeon Ye, Seungone Kim, Vipul

Gupta, Vivek Sharma, Xuhui Zhou, Caiming Xiong, Luis Villa, Stella Biderman, Alex Pentland, Sara Hooker, and Jad Kabbara. Bridging the data provenance gap across text, speech and video, 2025. URL https://arxiv.org/abs/2412.17847.

Holy Lovenia, Rahmad Mahendra, Salsabil Maulana Akbar, Lester James Validad Miranda, Jennifer Santoso, Elyanah Aco, Akhdan Fadhilah, Jonibek Mansurov, Joseph Marvin Imperial, Onno P. Kampman, Joel Ruben Antony Moniz, Muhammad Ravi Shulthan Habibi, Frederikus Hudi, Jann Railey Montalan, Ryan Ignatius Hadiwijaya, Joanito Agili Lopo, William Nixon, Börje F. Karlsson, James Jaya, Ryandito Diandaru, Yuze Gao, Patrick Amadeus Irawan, Bin Wang, Jan Christian Blaise Cruz, Chenxi Whitehouse, Ivan Halim Parmonangan, Maria Khelli, Wenyu Zhang, Lucky Susanto, Reynard Adha Ryanda, Sonny Lazuardi Hermawan, Dan John Velasco, Muhammad Dehan Al Kautsar, Willy Fitra Hendria, Yasmin Moslem, Noah Flynn, Muhammad Farid Adilazuarda, Haochen Li, Johanes Lee, R. Damanhuri, Shuo Sun, Muhammad Reza Qorib, Amirbek Djanibekov, Wei Qi Leong, Quyet V. Do, Niklas Muennighoff, Tanrada Pansuwan, Ilham Firdausi Putra, Yan Xu, Tai Ngee Chia, Ayu Purwarianti, Sebastian Ruder, William Chandra Tjhi, Peerat Limkonchotiwat, Alham Fikri Aji, Sedrick Keh, Genta Indra Winata, Ruochen Zhang, Fajri Koto, Zheng Xin Yong, and Samuel Cahyawijaya. SEACrowd: A multilingual multimodal data hub and benchmark suite for Southeast Asian languages. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 5155–5203, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.296. URL https://aclanthology.org/2024.emnlp-main.296/.

Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 2507–2521. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/11332b6b6cf4485b84afadb1352d3a9a-Paper-Conference.pdf.

Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023.

Alexandra Luccioni and Joseph Viviano. What's in the box? an analysis of undesirable content in the Common Crawl corpus. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pp. 182–189, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-short.24. URL https://aclanthology.org/2021.acl-short.24/.

Li Lucy, Suchin Gururangan, Luca Soldaini, Emma Strubell, David Bamman, Lauren Klein, and Jesse Dodge. AboutMe: Using self-descriptions in webpages to document the effects of English pretraining data filters. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 7393–7420, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.400. URL https://aclanthology.org/2024.acl-long.400/.

Muhammad Maaz, Hanoona Abdul Rasheed, Abdelrahman M. Shaker, Salman H. Khan, Hisham Cholakkal, Rao Muhammad Anwer, Timothy Baldwin, Michael Felsberg, and Fahad Shahbaz

Khan. Palo: A polyglot large multimodal model for 5b people. *ArXiv*, abs/2402.14818, 2024. URL https://api.semanticscholar.org/CorpusID:267782854.

Shravan Nayak, Kanishk Jain, Rabiul Awal, Siva Reddy, Sjoerd Van Steenkiste, Lisa Anne Hendricks, Karolina Stanczak, and Aishwarya Agrawal. Benchmarking vision language models for cultural understanding. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 5769–5790, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.329. URL https://aclanthology.org/2024.emnlp-main.329/.

Wilhelmina Nekoto, Vukosi Marivate, Tshinondiwa Matsila, Timi Fasubaa, Taiwo Fagbohungbe, Solomon Oluwole Akinola, Shamsuddeen Muhammad, Salomon Kabongo Kabenamualu, Salomey Osei, Freshia Sackey, Rubungo Andre Niyongabo, Ricky Macharm, Perez Ogayo, Orevaoghene Ahia, Musie Meressa Berhe, Mofetoluwa Adeyemi, Masabata Mokgesi-Selinga, Lawrence Okegbemi, Laura Martinus, Kolawole Tajudeen, Kevin Degila, Kelechi Ogueji, Kathleen Siminyu, Julia Kreutzer, Jason Webster, Jamiil Toure Ali, Jade Abbott, Iroro Orife, Ignatius Ezeani, Idris Abdulkadir Dangana, Herman Kamper, Hady Elsahar, Goodness Duru, Ghollah Kioko, Murhabazi Espoir, Elan van Biljon, Daniel Whitenack, Christopher Onyefuluchi, Chris Chinenye Emezue, Bonaventure F. P. Dossou, Blessing Sibanda, Blessing Bassey, Ayodele Olabiyi, Arshath Ramkilowan, Alp Öktem, Adewale Akinfaderin, and Abdallah Bashir. Participatory research for low-resourced machine translation: A case study in African languages. In Trevor Cohn, Yulan He, and Yang Liu (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 2144–2160, Online, November 2020. Association for Computational Linguistics. doi: 10.18653 /v1/2020.findings-emnlp.195. URL https://aclanthology.org/2020.findings-emnlp.195/.

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer,

Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. URL https://arxiv.org/abs/2303.08774.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, ACL '02, pp. 311–318, USA, 2002. Association for Computational Linguistics. doi: 10.3115/1073083.1073135. URL https://doi.org/10.3115/1073083.1073135.

Wannaphong Phatthiyaphaibun, Surapon Nonesung, Patomporn Payoungkhamdee, Peerat Limkonchotiwat, Can Udomcharoenchaikit, Jitkapat Sawatphol, Chompakorn Chaksangchaichot, Ekapol Chuangsuwanich, and Sarana Nutanong. Wangchanlion and wangchanx mrc eval, 2024. URL https://arxiv.org/abs/2403.16127.

Luiza Pozzobon, Patrick Lewis, Sara Hooker, and Beyza Ermis. From one to many: Expanding the scope of toxicity mitigation in language models, 2024. URL https://arxiv.org/abs/2403.03893.

Chen Qiu, Dan Oneață, Emanuele Bugliarello, Stella Frank, and Desmond Elliott. Multilingual multimodal learning with machine translated text. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp. 4178–4193, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-emnlp.308. URL https://aclanthology.org/2022.findings-emnlp.308/.

Qwen-Team. Qwen2.5-vl, January 2025. URL https://qwenlm.github.io/blog/qwen2.5-vl/.

Rita Ramos, Emanuele Bugliarello, Bruno Martins, and Desmond Elliott. PAELLA: Parameter-efficient lightweight language-agnostic captioning model. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Findings of the Association for Computational Linguistics: NAACL 2024*,

pp. 3549–3564, Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-naacl.225. URL https://aclanthology.org/2024.findings-naacl.225/.

Angelika Romanou, Negar Foroutan, Anna Sotnikova, Zeming Chen, Sree Harsha Nelaturu, Shivalika Singh, Rishabh Maheshwary, Micol Altomare, Mohamed A. Haggag, Snegha A, Alfonso Amayuelas, Azril Hafizi Amirudin, Viraat Aryabumi, Danylo Boiko, Michael Chang, Jenny Chim, Gal Cohen, Aditya Kumar Dalmia, Abraham Diress, Sharad Duwal, Daniil Dzenhaliou, Daniel Fernando Erazo Florez, Fabian Farestam, Joseph Marvin Imperial, Shayekh Bin Islam, Perttu Isotalo, Maral Jabbarishiviari, Börje F. Karlsson, Eldar Khalilov, Christopher Klamm, Fajri Koto, Dominik Krzemiński, Gabriel Adriano de Melo, Syrielle Montariol, Yiyang Nan, Joel Niklaus, Jekaterina Novikova, Johan Samir Obando Ceron, Debjit Paul, Esther Ploeger, Jebish Purbey, Swati Rajwal, Selvan Sunitha Ravi, Sara Rydell, Roshan Santhosh, Drishti Sharma, Marjana Prifti Skenduli, Arshia Soltani Moakhar, Bardia Soltani Moakhar, Ran Tamir, Ayush Kumar Tarun, Azmine Toushik Wasi, Thenuka Ovin Weerasinghe, Serhan Yilmaz, Mike Zhang, Imanol Schlag, Marzieh Fadaee, Sara Hooker, and Antoine Bosselut. Include: Evaluating multilingual language understanding with regional knowledge. *arXiv preprint arXiv:2411.19799*, 2024. URL https://arxiv.org/abs/2411.19799.

David Romero, Chenyang Lyu, Haryo Akbarianto Wibowo, Teresa Lynn, Injy Hamed, Aditya Nanda Kishore, Aishik Mandal, Alina Dragonetti, Artem Abzaliev, Atnafu Lambebo Tonja, Bontu Fufa Balcha, Chenxi Whitehouse, Christian Salamea, Dan John Velasco, David Ifeoluwa Adelani, David Le Meur, Emilio Villa-Cueva, Fajri Koto, Fauzan Farooqui, Frederico Belcavello, Ganzorig Batnasan, Gisela Vallejo, Grainne Caulfield, Guido Ivetta, Haiyue Song, Henok Biadglign Ademtew, Hernán Maina, Holy Lovenia, Israel Abebe Azime, Jan Christian Blaise Cruz, Jay Gala, Jiahui Geng, Jesus-German Ortiz-Barajas, Jinheon Baek, Jocelyn Dunstan, Laura Alonso Alemany, Kumaranage Ravindu Yasas Nagasinghe, Luciana Benotti, Luis Fernando D'Haro, Marcelo Viridiano, Marcos Estecha-Garitagoitia, Maria Camila Buitrago Cabrera, Mario Rodríguez-Cantelar, Mélanie Jouitteau, Mihail Mihaylov, Mohamed Fazli Mohamed Imam, Muhammad Farid Adilazuarda, Munkhjargal Gochoo, Munkh-Erdene Otgonbold, Naome Etori, Olivier Niyomugisha, Paula Mónica Silva, Pranjal Chitale, Raj Dabre, Rendi Chevi, Ruochen Zhang, Ryandito Diandaru, Samuel Cahyawijaya, Santiago Góngora, Soyeong Jeong, Sukannya Purkayastha, Tatsuki Kuribayashi, Thanmay Jayakumar, Tiago Timponi Torrent, Toqeer Ehsan, Vladimir Araujo, Yova Kementchedjhieva, Zara Burzo, Zheng Wei Lim, Zheng Xin Yong, Oana Ignat, Joan Nwatu, Rada Mihalcea, Thamar Solorio, and Alham Fikri Aji. Cvqa: Culturally-diverse multilingual visual question answering benchmark. *NeurIPS 2024*, 2024.

Florian Schneider and Sunayana Sitaram. M5 – a diverse benchmark to assess the performance of large multimodal models across multilingual and multicultural vision-language tasks. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 4309–4345, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.250. URL https://aclanthology.org/2024.findings-emnlp.250/.

Florian Schneider, Carolin Holtermann, Chris Biemann, and Anne Lauscher. Gimmick – globally inclusive multimodal multitask cultural knowledge benchmarking, 2025. URL https://arxiv.org/abs/2502.13766.

Shivalika Singh, Angelika Romanou, Clémentine Fourrier, David I Adelani, Jian Gang Ngui, Daniel Vila-Suero, Peerat Limkonchotiwat, Kelly Marchisio, Wei Qi Leong, Yosephine Susanto, et al.

Global mmlu: Understanding and addressing cultural and linguistic biases in multilingual evaluation. *arXiv preprint arXiv:2412.03304*, 2024a.

Shivalika Singh, Freddie Vargus, Daniel D'souza, Börje Karlsson, Abinaya Mahendiran, Wei-Yin Ko, Herumb Shandilya, Jay Patel, Deividas Mataciunas, Laura O'Mahony, Mike Zhang, Ramith Hettiarachchi, Joseph Wilson, Marina Machado, Luisa Moura, Dominik Krzemiński, Hakimeh Fadaei, Irem Ergun, Ifeoma Okoh, Aisha Alaagib, Oshan Mudannayake, Zaid Alyafeai, Vu Chien, Sebastian Ruder, Surya Guthikonda, Emad Alghamdi, Sebastian Gehrmann, Niklas Muennighoff, Max Bartolo, Julia Kreutzer, Ahmet Üstün, Marzieh Fadaee, and Sara Hooker. Aya dataset: An open-access collection for multilingual instruction tuning. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 11521–11567, Bangkok, Thailand, August 2024b. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.620. URL https://aclanthology.org/2024.acl-long.620/.

Shivalika Singh, Angelika Romanou, Clémentine Fourrier, David I. Adelani, Jian Gang Ngui, Daniel Vila-Suero, Peerat Limkonchotiwat, Kelly Marchisio, Wei Qi Leong, Yosephine Susanto, Raymond Ng, Shayne Longpre, Wei-Yin Ko, Sebastian Ruder, Madeline Smith, Antoine Bosselut, Alice Oh, Andre F. T. Martins, Leshem Choshen, Daphne Ippolito, Enzo Ferrante, Marzieh Fadaee, Beyza Ermis, and Sara Hooker. Global mmlu: Understanding and addressing cultural and linguistic biases in multilingual evaluation, 2025. URL https://arxiv.org/abs/2412.03304.

Ray Smith. An overview of the Tesseract OCR engine. In *Proc. Ninth Int. Conference on Document Analysis and Recognition (ICDAR)*, pp. 629–633, 2007.

Ashish V. Thapliyal, Jordi Pont Tuset, Xi Chen, and Radu Soricut. Crossmodal-3600: A massively multilingual multimodal evaluation dataset. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 715–729, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.45. URL https://aclanthology.org/2022.emnlp-main.45/.

Jesse Thomason, Daniel Gordon, and Yonatan Bisk. Shifting the baseline: Single modality performance on visual navigation & QA. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 1977–1983, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1197. URL https://aclanthology.org/N19-1197/.

Ahmet Üstün, Viraat Aryabumi, Zheng Yong, Wei-Yin Ko, Daniel D'souza, Gbemileke Onilude, Neel Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, Freddie Vargus, Phil Blunsom, Shayne Longpre, Niklas Muennighoff, Marzieh Fadaee, Julia Kreutzer, and Sara Hooker. Aya model: An instruction finetuned open-access multilingual language model. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15894–15939, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.845. URL https://aclanthology.org/2024.acl-long.845/.

Özlem Uzuner, Imre Solti, Fei Xia, and Eithon Cadag. Community annotation experiment for ground truth generation for the i2b2 medication challenge. *Journal of the American Medical*

*Informatics Association*, 17(5):519–523, 09 2010. ISSN 1067-5027. doi: 10.1136/jamia.2010.004 200. URL https://doi.org/10.1136/jamia.2010.004200.

Emiel van Miltenburg, Desmond Elliott, and Piek Vossen. Cross-linguistic differences and similarities in image descriptions. In Jose M. Alonso, Alberto Bugarín, and Ehud Reiter (eds.), *Proceedings of the 10th International Conference on Natural Language Generation*, pp. 21–30, Santiago de Compostela, Spain, September 2017. Association for Computational Linguistics. doi: 10.18653/v 1/W17-3503. URL https://aclanthology.org/W17-3503/.

Eva Vanmassenhove, Dimitar Shterionov, and Matthew Gwilliam. Machine translationese: Effects of algorithmic bias on linguistic complexity in machine translation. In Paola Merlo, Jorg Tiedemann, and Reut Tsarfaty (eds.), *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pp. 2203–2213, Online, April 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.eacl-main.188. URL https://aclanthology.org/2021.eacl-main.188/.

Ashmal Vayani, Dinura Dissanayake, Hasindri Watawana, Noor Ahsan, Nevasini Sasikumar, Omkar Thawakar, Henok Biadglign Ademtew, Yahya Hmaiti, Amandeep Kumar, Kartik Kuckreja, Mykola Maslych, Wafa Al Ghallabi, Mihail Mihaylov, Chao Qin, Abdelrahman M. Shaker, Mike Zhang, Mahardika Krisna Ihsani, Amiel Esplana, Monil Gokani, Shachar Mirkin, Harsh Singh, Ashay Srivastava, Endre Hamerlik, Fathinah Asma Izzati, Fadillah Adamsyah Maani, Sebastian Cavada, Jenny Chim, Rohit Gupta, Sanjay Manjunath, Kamila Zhumakhanova, Feno Heriniaina Rabevohitra, Azril Hafizi Amirudin, Muhammad Ridzuan, Daniya Najiha Abdul Kareem, Ketan More, Kunyang Li, Pramesh Shakya, Muhammad Saad, Amirpouya Ghasemaghaei, Amirbek Djanibekov, Dilshod Azizov, Branislava Jankovic, Naman Bhatia, Alvaro Cabrera, Johan Obando-Ceron, Olympiah Otieno, Fabian Farestam, Muztoba Rabbani, Sanoojan Baliah, Santosh Sanjeev, Abduragim Shtanchaev, Maheen Fatima, Thao Nguyen, Amrin Kareem, Toluwani Aremu, Nathan Xavier, Amit Bhatkal, Hawau Olamide Toyin, Aman Chadha, Hisham Cholakkal, Rao Muhammad Anwer, Michael Felsberg, Jorma Laaksonen, Thamar Solorio, Monojit Choudhury, Ivan Laptev, Mubarak Shah, Salman Khan, and Fahad Shahbaz Khan. All languages matter: Evaluating lmms on culturally diverse 100 languages. *ArXiv*, abs/2411.16508, 2024. URL https://api.semanticscholar.org/CorpusID:274234962.

Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4566–4575, 2015.

Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Mingjie Zhan, and Hongsheng Li. Measuring multimodal mathematical reasoning with math-vision dataset, 2024a.

Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, Tianle Li, Max KU, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhu Chen. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 95266–95290. Curran Associates, Inc., 2024b. URL https://proceedings.neurips.cc/pap er_files/paper/2024/file/ad236edc564f3e3156e1b2feafb99a24-Paper-Datasets_and_Be nchmarks_Track.pdf.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language

models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2022. Curran Associates Inc. ISBN 9781713871088.

Peng Xu, Wenqi Shao, Kaipeng Zhang, Peng Gao, Shuo Liu, Meng Lei, Fanqing Meng, Siyuan Huang, Yu Qiao, and Ping Luo. Lvlm-ehub: A comprehensive evaluation benchmark for large vision-language models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 47(3): 1877–1893, 2025. doi: 10.1109/TPAMI.2024.3507000.

Weihao Xuan, Rui Yang, Heli Qi, Qingcheng Zeng, Yunze Xiao, Yun Xing, Junjue Wang, Huitao Li, Xin Li, Kunyu Yu, et al. Mmlu-prox: A multilingual benchmark for advanced large language model evaluation. *arXiv preprint arXiv:2503.10497*, 2025.

Xiang Yue, Yuansheng Ni, Tianyu Zheng, Kai Zhang, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9556–9567, 2024a. doi: 10.1109/CVPR52733.2024.00913.

Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang, Kai Zhang, Shengbang Tong, Yuxuan Sun, Botao Yu, Ge Zhang, Huan Sun, Yu Su, Wenhu Chen, and Graham Neubig. Mmmu-pro: A more robust multi-discipline multimodal understanding benchmark, 2024b. URL https://arxiv.org/abs/2409.02813.

Xiang Yue, Yueqi Song, Akari Asai, Seungone Kim, Jean de Dieu Nyandwi, Simran Khanuja, Anjali Kantharuban, Lintang Sutawika, Sathyanarayanan Ramamoorthy, and Graham Neubig. Pangea: A fully open multilingual multimodal llm for 39 languages, 2025. URL https://arxiv.org/abs/2410.16153.

Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Y. Zou. When and why vision-language models behave like bags-of-words, and what to do about it? *ArXiv*, abs/2210.01936, 2022. URL https://api.semanticscholar.org/CorpusID:252734947.

Yuhang Zang, Wei Li, Jun Han, Kaiyang Zhou, and Chen Change Loy. Contextual object detection with multimodal large language models. *International Journal of Computer Vision*, pp. 1–19, 2024.

Mike Zhang and Antonio Toral. The effect of translationese in machine translation test sets. In Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, André Martins, Christof Monz, Matteo Negri, Aurélie Névéol, Mariana Neves, Matt Post, Marco Turchi, and Karin Verspoor (eds.), *Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)*, pp. 73–81, Florence, Italy, August 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-5208. URL https://aclanthology.org/W19-5208/.

Wenxuan Zhang, Sharifah Mahani Aljunied, Chang Gao, Yew Ken Chia, and Lidong Bing. M3exam: A multilingual, multimodal, multilevel benchmark for examining large language models. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023*, New Orleans, LA, USA, 2023. URL https://arxiv.org/abs/2306.05179.

Mingyang Zhou, Luowei Zhou, Shuohang Wang, Yu Cheng, Linjie Li, Zhou Yu, and Jingjing Liu. Uc2: Universal cross-lingual cross-modal vision-and-language pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4155–4165, June 2021.

# A   Data Collection Details

## A.1   Dataset Fields

| Field | Description |
|---|---|
| language | The language in which the question is written (e.g., "en" for English). |
| country | The country where the exam originated (e.g., "United States"). |
| contributor_country | The contributor's country of residence (e.g., "Spain"). |
| file_name | The internal database filename for the original exam document. |
| source | The URL or reference to the original exam document. |
| license | Licensing information of the exam (e.g., "Unknown" if not stated). |
| level | The educational level of the exam (e.g., "University Entrance"). |
| category_en | The exam subject category in English (e.g., "Chemistry"). |
| category_source_lang | The subject category as written in the original language (language). |
| original_question_num | The original question number in the source document. |
| question | The text of the question. |
| options | A list of possible answer choices. For example, ["Option A", "Option B", "Option C", "Option D"]. |
| answer | The index of the correct answer (e.g., 3 for the fourth option). |
| question_image | The extracted diagram, graph, or table associated with the question. |
| image_information | A label indicating the importance of the question_image for answering the question. Possible values include:<br>• "useful" - The image provides additional clarification.<br>• "essential" - The image is necessary to answer the question. |
| image_type | The category of question_image (e.g., "figure", "graph", "table") as described in Appendix A.2. |

Table 7: Structured dataset fields with descriptions used in data collection protocol.

## A.2   Categories of Visual Elements

We group the visual elements into eight primary categories in KALEIDOSCOPE. If an image falls into multiple categories, we assign the most representative based on the image's content.

| Visual Element Category | Question Image | Question and Answer |
|---|---|---|
| **Diagram.** Technical or schematic drawings illustrating processes, structures, or concepts. |  | **Question:** Wie verhält sich die Verarmungszone in der hier dargestellten Halbleiterdiode? <br> **Options:** <br> **A. Sie erweitert sich.** <br> B. Sie verengt sich. <br> C. Sie verändert sich nicht. <br> D. Sie verschwindet. |
| **Figure.** Illustrations, drawings, or visual representations of objects, patterns, or symbols. |  | **Question:** Applicable for D of stem 'B'- <br> **Options:** <br> A. contains more genes <br> B. unable to replicate <br> C. present in the nucleus <br> **D. used as a vector** |
| **Charts.** Images showing data plotted on axes, such as line graphs, bar charts, scatter plots, pie charts, flowcharts, organizational charts, and so on. |  | **Question:** Em uma xícara que já contém certa quantidade de açúcar, despeja-se café. A curva abaixo representa a função exponencial M(t), que fornece a quantidade de açúcar não dissolvido (em gramas), t minutos após o café ser despejado. Pelo gráfico, podemos concluir que. <br> **Options:** <br> **A.** $m(t) = 2^{(4-t/75)}$ <br> B. $m(t) = 2^{(4-t/50)}$ <br> C. $m(t) = 2^{(5-t/50)}$ <br> D. $m(t) = 2^{(5-t/150)}$ |
| **Map.** Geographical or spatial representations. |  | **Question:** Діяльність якого гетьмана можна характеризувати, спираючись на подану карту? <br> **Options:** <br> А. Б. Хмельницького <br> **В. І. Виговського** <br> С. Д. Многогрішного <br> D. І. Самойловича |

| Visual Element Category | Question Image | Question and Answer |
|---|---|---|
| **Photographs.** Photographic images of real-world scenes, objects, or people. |  | **Question:** Wat kun je zeggen over het verzorgingsgebied van deze McDonald's in Arnhem?<br>**Options:**<br>**A. Het verzorgingsgebied beperkt zich tot de stad Arnhem.**<br>B. Het verzorgingsgebied beperkt zich tot de provincie Gelderland.<br>C. Het verzorgingsgebied beperkt zich tot Nederland.<br>D. Het verzorgingsgebied beperkt zich tot de regio Arnhem en omstreken. |
| **Formula.** Mathematical equations, chemical formulas, mathematical diagrams, or related concepts. | $$2HI(g) \rightleftharpoons H_2(g) + I_2(g)$$ | **Question:** अभिक्रिया इस छवि में दिखाए गए समीकरण की विघटन की कोटि, साम्यावस्था स्थिरांक $K_p$ में सम्बद्ध है।<br>**Options:**<br>A. $\sqrt{\frac{1+2K_p}{2}}$<br>B. $\frac{1+2K_p}{2}$<br>C. $\frac{2K_p}{1+2K_p}$<br>**D. $\frac{2\sqrt{K_p}}{1+2\sqrt{K_p}}$** |
| **Table.** Structured data arranged in rows and columns. |  | **Question:**<br>به چند طریق می‌توان جدول نیم‌پر روبه‌رو را با عددهای ۱ تا ۴ طوری پر کرد که در هیچ سطر و ستونی عدد تکراری نداشته باشیم؟<br>**Options:**<br>A. 0<br>B. 1<br>**C. 2**<br>D. !4 |
| **Text.** Images containing primarily textual information. | ABC ত্রিভুজে B কোণের পরিমাণ ৪৮° এবং AB=AC। | **Question:** যদি E এবং F AB এবং AC-কে এমনভাবে ছেদ করে যেন EF ‖ BC হয়, তাহলে<br>**Options:**<br>**A. ১৩২°**<br>B. ১৬০°<br>C. ১৮০°<br>D. ১০৮° |

Table 8: Types of visual elements or images in the KALEIDOSCOPE benchmark. The correct answer is highlighted in **Bold Green**. Some samples are reformatted for better presentation.

## A.3 Complete Results

We report full multimodal performances in table 9 for each language and in table 10 for each subject. Each table reports, for each model and category; **Total Accuracy %**: the accuracy over all samples, **Valid Accuracy %**: the accuracy over successfully extracted answers and **Format Error % (FE)**: the proportion of unextracted answers.

| | | Latin Script | | | | | | Non-Latin Script | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | English | French | German | Dutch | Portuguese | Spanish | Arabic | Bengali | Croatian | Hindi | Hungarian | Lithuanian | Nepali | Persian | Russian | Serbian | Telugu | Ukrainian |
| **Gemini 1.5 Pro** | Total Acc. | 62.7 | 54.6 | 52.6 | 61.5 | 81.8 | 78.5 | 44.5 | 51.8 | 46.9 | 62.6 | 39.1 | 75.0 | 22.2 | 41.2 | 45.0 | 41.9 | 58.1 | 70.3 |
| | Valid Acc. | 63.2 | 55.2 | 52.6 | 62.5 | 83.4 | 78.5 | 44.7 | 52.7 | 48.1 | 63.2 | 40.5 | 75.0 | 22.8 | 42.1 | 46.2 | 43.6 | 58.3 | 70.3 |
| | FE | 0.9 | 1.0 | 0.0 | 1.6 | 1.9 | 0.0 | 0.5 | 1.8 | 2.5 | 0.9 | 3.4 | 0.0 | 2.4 | 2.1 | 2.6 | 3.8 | 0.4 | 0.0 |
| **Claude 3.5 Sonnet** | Total Acc. | 36.9 | 51.7 | 73.7 | 65.2 | 83.8 | 77.6 | 50.3 | 49.5 | 46.9 | 57.1 | 38.8 | 81.2 | 27.8 | 45.6 | 45.3 | 38.9 | 56.0 | 75.2 |
| | Valid Acc. | 63.2 | 51.7 | 73.7 | 65.6 | 83.8 | 77.7 | 50.3 | 49.6 | 47.5 | 57.2 | 38.9 | 81.4 | 28.0 | 45.6 | 45.4 | 39.6 | 56.0 | 75.2 |
| | FE | 41.6 | 0.0 | 0.0 | 0.6 | 0.0 | 0.1 | 0.0 | 0.2 | 1.2 | 0.1 | 0.4 | 0.3 | 0.8 | 0.0 | 0.2 | 1.8 | 0.0 | 0.0 |
| **GPT-4o** | Total Acc. | 42.6 | 46.2 | 52.4 | 60.1 | 76.6 | 73.8 | 41.9 | 60.2 | 36.4 | 53.0 | 36.4 | 76.2 | 19.0 | 41.3 | 37.6 | 32.4 | 41.6 | 68.5 |
| | Valid Acc. | 64.5 | 49.2 | 53.7 | 66.8 | 81.0 | 78.6 | 49.7 | 62.8 | 40.4 | 59.0 | 40.7 | 79.7 | 20.5 | 42.6 | 40.3 | 38.7 | 47.9 | 77.5 |
| | FE | 33.9 | 6.0 | 2.5 | 10.0 | 5.4 | 6.1 | 15.7 | 4.0 | 9.9 | 10.1 | 10.5 | 4.4 | 7.1 | 3.0 | 6.7 | 16.2 | 13.2 | 11.6 |
| **Qwen2.5-VL-72B** | Total Acc. | 53.8 | 46.2 | 49.3 | 57.4 | 73.3 | 72.5 | 36.1 | 49.5 | 33.3 | 46.2 | 37.5 | 69.1 | 23.8 | 35.4 | 39.3 | 35.9 | 47.0 | 65.2 |
| | Valid Acc. | 53.8 | 46.4 | 49.3 | 57.5 | 73.3 | 72.5 | 36.1 | 49.5 | 33.3 | 46.2 | 37.6 | 69.1 | 23.8 | 35.4 | 39.3 | 35.9 | 47.0 | 65.2 |
| | FE | 0.0 | 0.5 | 0.0 | 0.2 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.2 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **Aya-Vision-32B** | Total Acc. | 32.9 | 27.6 | 39.1 | 47.5 | 57.4 | 53.2 | 30.4 | 26.0 | 29.6 | 32.9 | 26.1 | 48.5 | 22.2 | 30.3 | 29.0 | 28.4 | 34.8 | 47.1 |
| | Valid Acc. | 33.0 | 29.3 | 39.5 | 48.7 | 58.8 | 55.5 | 30.5 | 26.1 | 29.6 | 33.7 | 26.5 | 48.5 | 22.2 | 31.1 | 29.4 | 28.5 | 34.8 | 47.1 |
| | FE | 0.4 | 6.0 | 1.1 | 2.4 | 2.4 | 4.2 | 0.5 | 0.2 | 0.0 | 2.3 | 1.6 | 0.0 | 0.0 | 2.5 | 1.3 | 0.4 | 0.1 | 0.0 |
| **Aya-Vision-8B** | Total Acc. | 28.9 | 30.2 | 32.4 | 42.2 | 47.6 | 46.3 | 27.2 | 18.2 | 29.0 | 29.2 | 27.1 | 34.4 | 23.8 | 27.7 | 25.2 | 28.5 | 11.1 | 36.6 |
| | Valid Acc. | 28.9 | 30.2 | 32.4 | 42.2 | 47.6 | 46.3 | 27.2 | 25.9 | 29.0 | 29.3 | 27.1 | 34.4 | 24.2 | 27.8 | 25.2 | 28.5 | 25.7 | 36.8 |
| | FE | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 29.5 | 0.0 | 0.2 | 0.0 | 0.0 | 1.6 | 0.3 | 0.0 | 0.0 | 56.8 | 0.6 |
| **Molmo-7B-D** | Total Acc. | 26.9 | 33.1 | 19.9 | 41.5 | 47.6 | 47.8 | 21.5 | 26.0 | 30.9 | 30.0 | 25.5 | 35.3 | 25.4 | 28.3 | 29.2 | 27.2 | 33.9 | 35.8 |
| | Valid Acc. | 27.0 | 33.3 | 19.9 | 41.5 | 47.6 | 47.8 | 21.5 | 26.0 | 30.9 | 30.1 | 25.5 | 35.3 | 25.4 | 28.3 | 29.2 | 27.2 | 33.9 | 35.8 |
| | FE | 0.2 | 0.8 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.2 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.1 |
| **Pangea-7B** | Total Acc. | 24.7 | 25.5 | 17.2 | 37.3 | 48.8 | 46.2 | 20.4 | 27.8 | 17.9 | 24.3 | 23.9 | 32.6 | 17.5 | 21.2 | 25.5 | 26.4 | 18.6 | 33.0 |
| | Valid Acc. | 27.9 | 31.1 | 19.6 | 39.0 | 51.0 | 49.4 | 22.3 | 31.8 | 21.0 | 29.7 | 26.2 | 33.4 | 23.9 | 25.2 | 28.9 | 30.1 | 33.9 | 33.9 |
| | FE | 11.4 | 18.1 | 12.2 | 4.3 | 4.3 | 6.5 | 8.4 | 12.8 | 14.8 | 18.3 | 8.8 | 2.4 | 27.0 | 16.0 | 11.9 | 12.4 | 45.2 | 2.6 |
| **Qwen2.5-VL-7B** | Total Acc. | 36.4 | 35.2 | 26.6 | 46.4 | 59.2 | 57.8 | 36.1 | 35.0 | 25.9 | 33.8 | 28.6 | 50.3 | 22.2 | 30.3 | 28.8 | 27.5 | 37.6 | 45.4 |
| | Valid Acc. | 36.4 | 35.3 | 26.6 | 46.4 | 59.2 | 57.8 | 36.1 | 35.1 | 26.1 | 33.9 | 28.6 | 50.3 | 22.2 | 30.3 | 28.8 | 27.5 | 37.6 | 45.4 |
| | FE | 0.0 | 0.3 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.6 | 0.2 | 0.2 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **Qwen2.5-VL-3B** | Total Acc. | 35.0 | 32.3 | 21.1 | 41.8 | 52.2 | 53.3 | 33.5 | 32.8 | 25.9 | 31.3 | 28.2 | 35.6 | 23.8 | 30.3 | 29.0 | 27.8 | 31.8 | 40.1 |
| | Valid Acc. | 35.0 | 32.4 | 21.1 | 41.8 | 52.3 | 53.3 | 33.7 | 32.8 | 26.4 | 31.3 | 29.2 | 35.6 | 23.8 | 30.4 | 29.0 | 28.1 | 31.8 | 40.1 |
| | FE | 0.0 | 0.3 | 0.0 | 0.0 | 0.1 | 0.0 | 0.5 | 0.2 | 1.9 | 0.0 | 3.2 | 0.0 | 0.0 | 0.2 | 0.0 | 1.1 | 0.0 | 0.0 |

Table 9: **Total Accuracy %**, **Valid Accuracy %** and **Format Error % (FE)** grouped by Language in KALEIDOSCOPE for multimodal samples.

| | | Biology | Chemistry | Driving License | Economics | Engineering | Geography | History | Language | Mathematics | Medicine | Physics | Reasoning | Social Sciences | Sociology |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Gemini 1.5 Pro** | Total Acc. | 60.1 | 60.2 | 64.4 | 64.1 | 57.0 | 72.8 | 78.7 | 83.5 | 46.9 | 69.6 | 57.4 | 51.2 | 85.7 | 92.3 |
| | Valid Acc. | 60.3 | 60.4 | 64.4 | 64.1 | 57.3 | 72.8 | 78.7 | 83.5 | 48.6 | 70.2 | 57.8 | 52.0 | 85.7 | 92.3 |
| | FE | 0.3 | 0.4 | 0.0 | 0.0 | 0.4 | 0.0 | 0.0 | 0.0 | 3.5 | 0.8 | 0.7 | 1.7 | 0.0 | 0.0 |
| **Claude 3.5 Sonnet** | Total Acc. | 61.6 | 59.0 | 64.2 | 63.4 | 50.0 | 81.4 | 83.7 | 85.1 | 43.7 | 72.9 | 58.4 | 52.2 | 82.9 | 91.0 |
| | Valid Acc. | 62.9 | 59.7 | 64.2 | 63.8 | 64.4 | 81.5 | 83.7 | 85.5 | 44.4 | 73.8 | 58.7 | 53.3 | 82.9 | 93.4 |
| | FE | 2.1 | 1.2 | 0.0 | 0.0 | 22.4 | 0.2 | 0.0 | 0.0 | 1.5 | 1.2 | 0.6 | 2.0 | 0.0 | 2.6 |
| **GPT-4o** | Total Acc. | 60.7 | 47.1 | 65.5 | 64.1 | 45.7 | 76.2 | 70.8 | 78.3 | 39.4 | 64.6 | 51.9 | 43.9 | 74.3 | 87.2 |
| | Valid Acc. | 63.9 | 52.9 | 73.1 | 66.7 | 56.3 | 80.4 | 86.4 | 85.8 | 44.0 | 75.6 | 54.7 | 51.0 | 88.1 | 93.2 |
| | FE | 5.1 | 11.1 | 10.4 | 3.8 | 18.9 | 5.2 | 18.1 | 8.7 | 10.5 | 14.6 | 5.1 | 13.9 | 15.7 | 6.4 |
| **Qwen2.5-VL-72B** | Total Acc. | 37.6 | 33.2 | 39.0 | 37.4 | 28.8 | 40.6 | 48.9 | 72.2 | 30.1 | 36.7 | 33.7 | 27.4 | 52.9 | 64.1 |
| | Valid Acc. | 37.6 | 33.2 | 39.0 | 37.7 | 28.9 | 40.7 | 48.9 | 72.2 | 30.4 | 36.7 | 33.7 | 27.4 | 52.9 | 64.1 |
| | FE | 0.0 | 0.0 | 0.0 | 0.8 | 0.2 | 0.2 | 0.0 | 0.0 | 0.9 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **Aya-Vision-32B** | Total Acc. | 40.2 | 34.6 | 47.1 | 29.8 | 34.6 | 50.5 | 61.4 | 71.0 | 28.8 | 52.1 | 31.8 | 27.5 | 64.3 | 70.5 |
| | Valid Acc. | 40.7 | 34.8 | 47.1 | 29.8 | 34.8 | 50.5 | 61.4 | 71.2 | 29.6 | 52.3 | 33.0 | 27.6 | 64.3 | 70.5 |
| | FE | 1.3 | 0.5 | 0.0 | 0.0 | 0.7 | 0.0 | 0.0 | 0.2 | 2.8 | 0.4 | 3.6 | 0.3 | 0.0 | 0.0 |
| **Aya-Vision-8B** | Total Acc. | 53.8 | 50.0 | 54.5 | 58.8 | 48.4 | 70.4 | 77.1 | 84.9 | 40.3 | 63.3 | 42.3 | 42.3 | 80.0 | 87.2 |
| | Valid Acc. | 53.8 | 50.0 | 54.5 | 58.8 | 48.4 | 70.4 | 77.1 | 85.1 | 40.3 | 63.3 | 42.3 | 42.3 | 80.0 | 87.2 |
| | FE | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.2 | 0.1 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **Molmo-7B-D** | Total Acc. | 35.2 | 33.5 | 44.9 | 27.5 | 24.4 | 37.6 | 42.1 | 60.1 | 28.7 | 40.4 | 26.7 | 27.5 | 51.4 | 60.3 |
| | Valid Acc. | 35.3 | 33.5 | 44.9 | 27.5 | 24.4 | 37.6 | 42.1 | 60.1 | 28.8 | 40.4 | 26.7 | 27.5 | 52.2 | 61.0 |
| | FE | 0.1 | 0.0 | 0.0 | 0.0 | 0.1 | 0.0 | 0.0 | 0.0 | 0.1 | 0.0 | 0.0 | 0.0 | 1.4 | 1.3 |
| **Pangea-7B** | Total Acc. | 30.9 | 22.0 | 38.2 | 32.1 | 21.4 | 35.6 | 42.1 | 65.8 | 24.8 | 35.0 | 24.7 | 21.9 | 50.0 | 55.1 |
| | Valid Acc. | 33.4 | 34.1 | 39.4 | 33.9 | 24.2 | 36.7 | 42.4 | 66.0 | 29.0 | 38.4 | 28.9 | 26.6 | 53.8 | 57.3 |
| | FE | 7.5 | 35.3 | 2.9 | 5.3 | 11.5 | 3.0 | 0.6 | 0.2 | 14.4 | 8.8 | 14.7 | 17.6 | 7.1 | 3.8 |
| **Qwen2.5-VL-7B** | Total Acc. | 34.3 | 16.2 | 41.2 | 22.1 | 30.3 | 38.7 | 45.3 | 56.6 | 28.6 | 35.4 | 27.1 | 24.3 | 57.1 | 57.7 |
| | Valid Acc. | 35.4 | 28.0 | 41.6 | 30.9 | 30.3 | 39.5 | 45.6 | 56.6 | 28.6 | 35.4 | 27.1 | 25.1 | 58.0 | 57.7 |
| | FE | 3.0 | 42.2 | 1.1 | 28.2 | 0.0 | 1.9 | 0.6 | 0.0 | 0.3 | 0.0 | 0.0 | 3.4 | 1.4 | 0.0 |
| **Qwen2.5-VL-3B** | Total Acc. | 42.6 | 38.5 | 44.9 | 42.7 | 32.4 | 51.0 | 52.9 | 75.7 | 30.1 | 45.8 | 34.7 | 29.5 | 68.6 | 73.1 |
| | Valid Acc. | 42.6 | 38.5 | 44.9 | 42.7 | 32.4 | 51.0 | 52.9 | 75.7 | 30.1 | 45.8 | 34.7 | 29.5 | 68.6 | 73.1 |
| | FE | 0.0 | 0.1 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.1 | 0.0 | 0.1 | 0.0 | 0.0 | 0.0 |

Table 10: **Total Accuracy %**, **Valid Accuracy %** and **Format Error % (FE)** grouped by Subject in KALEIDOSCOPE for multimodal samples.

We report full multimodal performances in table 9 for each language and in table 10 for each subject. Each table reports, for each model and category; **Total Accuracy %**: the accuracy over all samples, **Valid Accuracy %**: the accuracy over successfully extracted answers and **Format Error % (FE)**: the proportion of unextracted answers.

| | Qwen2.5-VL-7B | | | Gemini 1.5 Pro | | |
|---|---|---|---|---|---|---|
| | *Total Acc.* | *Valid Acc.* | *FR* | *Total Acc.* | *Valid Acc.* | *FR* |
| Diagram | 37.8 | 37.9 | 0.3 | 58.9 | 59.6 | 1.2 |
| Figure | 34.6 | 34.8 | 0.7 | 49.0 | 50.0 | 1.9 |
| Graph | 45.2 | 45.2 | 0.0 | 67.4 | 68.2 | 1.2 |
| Map | 46.7 | 46.7 | 0.0 | 70.9 | 70.9 | 0.0 |
| Photo | 53.9 | 54.1 | 0.5 | 74.2 | 74.3 | 0.2 |
| Formula | 37.0 | 37.3 | 0.8 | 66.7 | 68.7 | 2.9 |
| Table | 34.5 | 34.6 | 0.2 | 74.7 | 76.1 | 1.8 |
| Text | 79.8 | 79.8 | 0.0 | 83.7 | 83.7 | 0.0 |

Table 11: **Total Accuracy %**, **Valid Accuracy %** and **Format Error % (FE)** grouped by Image Type in KALEIDOSCOPE for captioning + OCR experiment.

## A.4  License

To ensure ethical data usage, we prioritize sources that permit redistribution and academic use. During data collection, we filter out content from sources with restrictive licensing policies. Addi-

tionally, our dataset does not include personally identifiable information, and all collected exams are either publicly available or obtained under appropriate agreements.

## A.5 Prompts

The prompts that we used to perform all experiments were designed to ensure consistency across languages. Examples are shown both in English and Spanish as an overview. Below is a summary of the key components.

### A.5.1 System Message

A system message sets the context for the model, instructing it to act as an expert in solving multiple-choice questions. For zero-shot CoT prompting, the message is provided in all the evaluation languages to support language-specific evaluation.

- **Zero-shot CoT**:

  - English: `You are an expert at solving multiple-choice questions. Carefully analyze the question, think step by step, and provide your FINAL answer between the tags <ANSWER> X </ANSWER>, where X is ONLY the correct choice. Do not write any additional text between the tags.`
  - Spanish: `Eres un experto en resolver preguntas de opción múltiple. Analiza cuidadosamente la pregunta, piensa paso a paso y proporciona tu respuesta FINAL entre las etiquetas <ANSWER> X </ANSWER>, donde X es ÚNICAMENTE la opción correcta. No escribas ningún texto adicional entre las etiquetas.`

- **Direct answer**:

  `You are a helpful assistant who answers multiple-choice questions. For each question, output your final answer in JSON format with the following structure: {"choice": "The correct option (e.g., A, B, C, or D)"}. ONLY output this format exactly. Do not include any additional text or explanations outside the JSON structure.`

### A.5.2 Keywords

Language-specific keywords are used to structure the prompts consistently across languages. These include terms for "Question," "Options," and "Answer" to be included when generating the prompt. For example:

- English: `{"question": "Question", "options": "Options", "answer": "Answer"}`

- Spanish: `{"question": "Pregunta", "options": "Opciones", "answer": "Respuesta"}`

### A.5.3 Prompt Examples

System messages A.5.1 and Keywords A.5.2 are used to systematically craft the prompt for a model in a specific language. We show examples of both a closed and an open model in Table 12.

| Open Model | Closed Model |
|---|---|
| | SYSTEM:<br>Eres un experto en resolver preguntas de opción múltiple. Analiza cuidadosamente la pregunta, piensa paso a paso y proporciona tu respuesta FINAL entre las etiquetas <ANSWER> X </ANSWER>, donde X es ÚNICAMENTE la opción correcta. No escribas ningún texto adicional entre las etiquetas.<br>USER: |
| SYSTEM:<br>You are a helpful assistant who answers multiple-choice questions. For each question, output your final answer in JSON format with the following structure: "choice":"The correct option (e.g., A, B, C, or D)". ONLY output this format exactly. Do not include any additional text or explanations outside the JSON structure. Output your choice in the specified JSON format.<br><br>USER:<br><br><br>Question: Make CORRECT match between Group-I and Group-II, in relation to interaction between two species.<br>Options:<br>A.) P-I, Q-II, R-III, S-IV<br>B.) P-III, Q-II, R-IV, S-I<br>C.) P-IV, Q-III, R-II, S-I<br>D.) P-III, Q-IV, R-II, S-I<br>Answer: | <br><br>Pregunta: Ante esta imagen en un paciente con un trastorno motor en miembros inferiores, señale la respuesta INCORRECTA:<br>Opciones:<br>A.) Debemos buscar una malformación de Chiari<br>B.) En algunos casos se asocia a hidrocefalia<br>C.) Se caracteriza por una pérdida de la sensibilidad táctil y vibratoria con preservación de la sensación térmica y dolorosa<br>D.) Puede producirse tras traumatismos o infecciones<br><br>Respuesta: |

Table 12: **Prompt examples in** KALEIDOSCOPE. Multimodal prompt samples with interleaved image are shown for an open model and a closed model.

## A.6 Captioning & OCR

We instantiated Gemini 1.5 Pro with the following instruction to generate synthetic captions from the images in KALEIDOSCOPE. Prompts with image augmentations are shown in Table 13.

| Open Model | Closed Model |
| --- | --- |

```
SYSTEM:
You are a helpful assistant who answers
multiple-choice questions.  For each
question, output your final answer in
JSON format with the following structure:
"choice":"The correct option (e.g., A,
B, C, or D)".  ONLY output this format
exactly.  Do not include any additional text
or explanations outside the JSON structure.


USER:
#include<stdio.h>

int main(int argc, char *argv[]){

    char a = 'P';

    char b = 'x';

    char c = (a & b) + '*';

    char d = (a | b) - '-';

    char e = (a ^ b) + '+';

    printf("%c %c %c\n", c, d, e);

    return 0;

}
ASCII encoding for relevant characters is given below
```

| A | B | C | ... | Z |
| --- | --- | --- | --- | --- |
| 65 | 66 | 67 | ... | 90 |

| a | b | c | ... | z |
| --- | --- | --- | --- | --- |
| 97 | 98 | 99 | ... | 122 |

| * | + | – |
| --- | --- | --- |
| 42 | 43 | 45 |

Caption: The code initializes character
variables 'a' to 'P' and 'b' to 'x'.
It then calculates 'c', 'd', and 'e'
using bitwise operations (&, |, ^) and
character addition with '*', '-', and '+',
respectively.  The `printf` function outputs
the characters c, d, and e.  Below the code,
three tables display ASCII values:  one for
uppercase letters 'A' to 'Z' (65 to 90),
another for lowercase letters 'a' to 'z'
(97 to 122), and a third for symbols '*',
'+', and '-' (42, 43, and 45, respectively).
Ellipses (...)  within the tables indicate
omitted values between the shown characters.
OCR: ##include<stdio.h>\}\n \nint main(int
argc,\n \nchar a = 'P';\nchar b = 'x';\nchar
c = (a &\nchar d = (a |\nchar e = (a \u
201c\n \nprintf (\"sc \%\nreturn 0;\n \n
\}\n \nchar *argv[]) \{\n \nby + te;\nb )
- '-\%3\nb ) + \"Hy\n se\\n\", c, d, e);\n
\nASCII encoding for relevant characters is
given below\n \n 42| 43) 45\n \n
Question:  What is printed by the following
ANSI C program? Options: A.) z K s B.) 122 75
83 C.) * - + D.) P x +
Answer:

SYSTEM:
Eres un experto en resolver preguntas de
opción múltiple.  Analiza cuidadosamente la
pregunta, piensa paso a paso y proporciona
tu respuesta FINAL entre las etiquetas
<ANSWER> X </ANSWER>, donde X es ÚNICAMENTE
la opción correcta.  No escribas ningún
texto adicional entre las etiquetas.


USER:
```

| Mes | Peso total |
| --- | --- |
| 1 | 1.500 gramos |
| 2 | 2.600 gramos |
| 3 | 3.700 gramos |
| 4 | 4.800 gramos |

```
Caption:  This table presents data on total
weight, measured in grams, across four
months.  The table consists of two columns:
Mes (Month) and Peso total (Total Weight).
Month 1 shows a weight of 1,500 grams, Month
2 shows 2,600 grams, Month 3 shows 3,700
grams, and Month 4 shows 4,800 grams.  The
table is a simple grid format with plain
black text on a white background.
OCR: Wes | Pesototal\n 1 [1500 gramos\n 2
|_2600 grmos\n 3 | 3700 gemoe\na \n \n \u
201cZOO grams\n \n
Pregunta:  Un perro cachorro tenía un peso
de 1.500 gramos al mes de nacido.  En la
tabla se muestra el peso del cachorro en
los primeros cuatro meses.  De acuerdo con
la tabla, ¿Cuál es el cambio del peso del
cachorro entre un mes y el mes siguiente?
Opciones:
A.) Disminuyó 1.500 gramos
B.) Disminuyó 2.600 gramos
C.) Aumentó 1.100 gramos
D.) Aumentó 3.300 gramos
Respuesta:
```

Table 13: **Caption+OCR prompt examples in** Kaleidoscope**.** Prompts are shown for open
and closed models in English and Spanish. Caption and OCR additions are highlighted in green.

**Gemini 1.5 Pro's prompt for captioning:**

```
**Instruction:**
You are an expert image captioner. Generate highly detailed, precise,
and academically relevant textual descriptions of images sourced from exam questions,
ensuring all critical visual elements are captured for accurate problem-solving.

**Guidelines:**

Exam-Specific Analysis:

- Primary Elements: Identify and describe key components (e.g., diagrams, charts, graphs,
labels, symbols, annotations) and their exact attributes (e.g., numerical values, units,
directional arrows, text annotations).

- Secondary Details: Note stylistic features (e.g., "black-and-white schematic,"
"color-coded bars in a graph"), spatial relationships (e.g., "force vectors
pointing northwest"), and contextual clues (e.g., axes labels, legends, scales).

- Textual Elements: Explicitly transcribe all visible text (e.g., labels like
"Mitochondria," numbers like "5V," titles like "Figure 2: Velocity vs. Time").

Academic Precision:

- Technical Focus: Prioritize details critical to exam questions (e.g., "a right
triangle with hypotenuse labeled c = 10 cm," "a bar graph comparing GDP of 5 countries,
with Japan's bar shaded blue at 4.3 trillion").

- Diagrams/Charts: Specify type (e.g., "pie chart," "circuit diagram") and components
(e.g., "resistor symbol connected to a battery").

- Scientific Relevance: Highlight measurements, units, symbols (e.g., "T = 25°C,"
"a pulley system with frictionless ropes").

Structure & Clarity:

- Begin with the image's purpose (e.g., "A biology diagram of a plant cell")
followed by a systematic breakdown (left-to-right, top-to-bottom, or by functional layers).

- Use neutral, objective language. Avoid assumptions unless implied by context (e.g.,
"a downward arrow labeled 9.8 m/s² likely representing gravitational acceleration").

**Output Format:**

- Single paragraph (4-6 sentences).

- Example:
```

"A physics diagram depicts two blocks on a frictionless inclined plane: Block A (5 kg)
is connected via a rope to Block B (3 kg) over a pulley. Angle theta = 30°, with
vectors labeled F_normal and F_gravity. A scale beside the plane shows time t = 0s
to t = 5s. Text at the bottom reads: 'Calculate tension in the rope.' The image is
monochrome, with dashed lines indicating motion direction."

Constraints:

- Avoid Omissions: Ensure no labels, numbers, or symbols are overlooked, even if
small or peripheral.

- Neutral Tone: Exclude subjective interpretations (e.g., "messy handwriting"
or "complex diagram") unless style is exam-relevant (e.g., "a hand-drawn sketch with
annotations").

## A.7   Open-Weight Models CoT Results

To benchmark the models, we initially designed a CoT prompt that instructed the models to think
step-by-step and then provide the correct answer, marking the choice with the tags `<ANSWER>`
`</ANSWER>`. However, in preliminary experiments, we found this instruction too complex for mid-
to small-sized models (32B–3B), which struggled to follow it consistently.

In Table 14, we compare results using CoT versus the direct English-language prompt adopted in
our final evaluation. The error rate was considerably higher for most models, even after cleaning and
extracting answers with regex matching their typical output formats. Two exceptions were Pangea
and Molmo, which showed lower error rates with the CoT prompt. However, this was because
both ignored the reasoning instruction and simply output the selected option, making extraction
easier. Prompt choice significantly impacted performance: the direct English prompt improved
results across all models except Pangea, which remained unchanged.

| | Overall CoT | | | Overall In-English | | |
| | | Valid Responses | | | Valid Responses | |
| **Model** | Acc. | F.E. | Valid Acc. | Acc. | F.E. | Valid Acc. |
|---|---|---|---|---|---|---|
| Aya-Vision-32B | 38.94 | 8.04 | 42.06 | 39.27 | 1.05 | 39.66 |
| Aya-Vision-8B | 33.08 | 6.22 | 35.15 | 35.09 | 0.07 | 35.11 |
| Molmo-7B-D | 32.86 | 0.01 | 32.87 | 32.87 | 0.04 | 32.88 |
| Pangea-7B | 31.24 | 5.61 | 33.45 | 31.31 | 7.42 | 34.02 |
| Qwen2.5-VL-7B | 35.18 | 6.34 | 37.64 | 39.56 | 0.08 | 39.60 |
| Qwen2.5-VL-3B | 32.90 | 1.40 | 33.33 | 35.56 | 0.19 | 35.63 |

Table 14: **Comparison of CoT and direct English prompting on** KALEIDOSCOPE **for small
models..** Reported values are macro-averaged accuracy (%) across all languages.

## A.8 Comparison with Other Benchmarks

Table A.8 offers a concise comparison of key multimodal benchmarks. MMMU (Yue et al., 2024a), SEED-Bench (Li et al., 2024a), and MME (Fu et al., 2023) are single-language datasets focused mainly on image-text pairs, with SEED-Bench also incorporating video-text. MME is notably smaller and only partially human-annotated, using mostly true/false formats. In contrast, M3Exam (Zhang et al., 2023), EXAMS-V (Das et al., 2024), and M5 (Schneider & Sitaram, 2024) introduce multilingualism—M5 being the most extensive with 41 languages—though much of its content is not multiple-choice and lacks verified annotations.

| Benchmark | Languages | Samples | Multimodal | Modalities | Human Annotation | Answer type |
|---|---|---|---|---|---|---|
| **MMMU** (Yue et al., 2024a) | 1 | 11,550 | 11,264 | Image-Text | Yes | MCQA |
| **SEED-Bench** (Li et al., 2024a) | 1 | 19,242 | 19,242 | Image-Text, Video-Text | Partial | MCQA |
| **MME** (Fu et al., 2023) | 1[†] | 2,194 | 0 | Image-Text | Partial | Y/N |
| **M3Exam** (Zhang et al., 2023) | 9 | 12,317 | 2,816 | Image-Text | Yes | MCQA |
| **EXAMS-V** (Das et al., 2024) | 11 | 20,932 | 5,086 | Image-Text | Yes | MCQA |
| **M5** (Schneider & Sitaram, 2024) | 41 | 237,094 | 1,422 | Image-Text | Yes | Mix |
| **KALEIDOSCOPE** | **18** | **20,911** | **11,459** | Image-Text | Yes | MCQA |

Table 15: **Comparison of Multimodal Benchmarks.** [†]All but 40 questions are in English that measure machine translation capability from Chinese to English.

KALEIDOSCOPE stands out by offering a balanced composition of 20,911 samples across 18 languages, with a strong focus on multimodal reasoning (11,459 Image-Text samples), comprehensive human annotation, and a consistent multiple-choice setup. Compared to existing benchmarks, KALEIDOSCOPE is more linguistically diverse than M3Exam and EXAMS-V, includes more multimodal samples than M5, and ensures higher quality through expert-verified annotations, making it a robust and equitable benchmark for evaluating multilingual multimodal models.

# B Evaluation Metrics and the MCQA Framework

Traditional evaluation metrics for VLMs, such as exact match accuracy, BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), and CIDEr (Vedantam et al., 2015), rely on surface-level n-gram comparisons that often penalize semantically equivalent answers phrased differently from reference texts. In contrast, the multiple-choice question answering (MCQA) framework (Hendrycks et al., 2021; Romero et al., 2024; Lu et al., 2022; Yue et al., 2024a) offers a more human-like evaluation paradigm by providing predefined answer options. This reduces ambiguity in scoring and facilitates the creation of evaluation datasets that capture both domain knowledge and linguistic/cultural nuances across languages. Although concerns regarding oversaturation and reliance on superficial cues in MCQA exist (Du et al., 2023; Yuksekgonul et al., 2022), these can be mitigated by extending the answer option space and applying rigorous filtering strategies (Wang et al., 2024b; Yue et al., 2024a). Our primary challenge lies in the scarcity of questions that are both multimodal and culturally agnostic. As demonstrated by results from KALEIDOSCOPE and related studies (Maaz et al., 2024;

[Nayak et al., 2024](#)), oversaturation is not a prevalent issue. Consequently, bridging this evaluation gap is of key importance. To ensure high data quality, source data in KALEIDOSCOPE are manually verified by qualified processors in accordance with established criteria ([2.2](#)), maintaining a clear distinction between verified and unverified data.

## B.1  Selected Dataset Samples

The subsequent table presents one sample from each dataset, including the question, the associated image, the provided answers, and the correct answer highlighted in green.

| Language | Question Image | Question and Answer |
|---|---|---|
| Arabic |  | **High School Exam**   **Physics**<br><br>**Question:**<br>قراءة الأميتر الحراري في الدائرة الموضحة تساوي ......<br>**Options:**<br>**A. 0.2 A**<br>B. 2 A<br>C. 0.02 A<br>D. 20 A |
| Arabic |  | **High School Exam**   **Geology**<br><br>**Question:**<br>من خلال الجدول استنتج الدولة الأقل تضرراً حال تعرضها لزلزال<br>**Options:**<br><br>A. ل<br><br>**B. ع**<br>C. ص<br>D. س |
| Bengali |  | **BRTA Driving Test**   **Driving**<br><br>**Question:** এই চিহ্নটি দ্বারা কি বুঝায়?<br>**Options:**<br>A. শুধুমাত্র সাইকেল চলাচলের জন্য<br>**B. সাইকেল চলাচল নিষেধ**<br>C. মোটরসাইকেল চলাচল নিষেধ<br>D. শুধুমাত্র মোটরসাইকেল চলাচলের জন্য |
| Bengali |  | **HSC Exam**   **Geography**<br><br>**Question:** উদ্দীপকের 'ক' ও 'খ' বায়ুপ্রবাহের সাধারণ বৈশিষ্ট্য- i. দক্ষিণ-পশ্চিম দিকে প্রবাহিত হয় ii. ডান দিকে বেঁকে যায় iii. সম উষ্ণতাবিশিষ্ট<br>নিচের কোনটি সঠিক?<br>**Options:**<br>A. i ও ii<br>B. i ও iii<br>**C. ii ও iii**<br>D. i, ii ও iii |
| Bengali |  | **BCS Exam**   **Reasoning**<br><br>**Question:** নিচের চিত্রে কয়টি ত্রিভুজ আছে?<br>**Options:**<br>A. ৫টি<br>B. ৬টি<br>**C. ৮টি**<br>D. ৪টি |

| Language | Question Image | Question and Answer |
|---|---|---|

**Dutch**

Dutch Central Exam · Economics

**Question:** Bekijk bovenstaand diagram. Hoeveel ton klein chemisch afval werd er in Nederland in 2000 ingezameld?
**Options:**
**A. "21.000 ton"**
B. "19.500 ton"
C. "23.000 ton"
D. "20.500 ton"

---

**English**

$$4x^2 - 9 = (px + t)(px - t)$$

SAT · Mathematics

**Question:** In the equation, p and t are constants. Which of the following could be the value of p?
**Options:**
**A. 2**
B. 3
C. 4
D. 9

---

**English**



UCEED Exam · Design

**Question:** Four spheres start revolving clockwise in concentric circles from their initial positions as shown below. Yellow travels at 2m/sec, green at 4m/sec, red at 2m/sec and blue at 4m/sec. Which of the following statement(s) is/are TRUE?
**Options:**
**A. Yellow and green never cross (overtake) each other**
B. Red and blue takes the same time to complete one revolution
C.Yellow takes less time than green to complete one revolution
D. Blue and red will cross each other twice after the first 3 complete revolutions of blue

---

**English**



HSC Exam · Biology

**Question:** Which type of food digest in lebelled 'S' mentioned in the figure?
**Options:**
A. Potato
**B. Pulse**
C. Oil
D. Ghee

Continued on next page

| Language | Question Image | Question and Answer |
|---|---|---|
| English |  | **GATE**　**Engineering**<br><br>**Question:** Consider the CMOS circuit shown in the figure (substrates are connected to their respective sources). The gate width $W$ to gate length $L$ ratios $W/L$ of the transistors are as shown. Both transistors have the same gate oxide capacitance per unit area. For the pMOSFET, the threshold voltage is $-1V$ and the mobility of holes is $40\,\text{cm}^2/\text{V.s}$. For the nMOSFET, the threshold voltage is $1V$ and the mobility of electrons is $300\,\text{cm}^2/\text{V.s}$. The steady-state output voltage $V_o$ is _____<br>**Options:**<br>A. equal to 0 V<br>B. more than 2 V<br>**C. less than 2 V**<br>D. equal to 2V |
| Flemish |  | **Physician Exam**　**Language**<br><br>**Question:** Figuur 1A toont de A-weging voor het menselijk gehoor. Deze figuur leert ons dat<br>**Options:**<br>**A. de mens tonen rond de 1000 Hz het beste hoort.**<br>B. mensen tonen van 10.000 Hz niet meer kunnen horen.<br>C. een toon met dezelfde fysische geluidssterkte altijd even intens wordt gehoord.<br>D. bij gelijke geluidssterkte, een mens 100 Hz zachter hoort dan 50 Hz. |
| French |  | **Mathematical Kangaroo**　**Mathematics**<br><br>**Question:** On attache ensemble des anneaux comme indiqué ci-contre de façon à former une chaîne de 1,7 m de longueur. Combien d'anneaux sont nécessaires ?<br>**Options:**<br>A. 30<br>**B. 42**<br>C. 21<br>D. 85 |
| German |  | **Amateur Radio Exam**　**Engineering**<br><br>**Question:** Wie groß ist die Gate-Source-Spannung, wenn sich der Schleifer von $R_3$ am Anschlag 1 befindet?<br>**Options:**<br>**A. 3,5 V**<br>B. 2,77 V<br>C. 3,7 V<br>D. 0,45 V |

| Language | Question Image | Question and Answer |
|---|---|---|

**Row 1**

Language: Hindi

Question Image:



Question and Answer:

[Science Olympiad] [Biology]

**Question:** अजय अपने घर से 20 मिनट पैदल चलकर दोपहर के 3.30 बजे सिनेमा हॉल पहुँचा। वह जल्दी-जल्दी सिनेमा हॉल में घुस गया। इसे अस-पास साफ-साफ देखने में कुछ समय लगा। इस दौरान उसकी आँखों में किस तरह के परिवर्तन आए होंगे?

**Options:**
A.  वृत्तीय और रेडियल मांसपेशियां शिथिल होती हैं जबकि पुतलियां संकुचित होती हैं।
**B.  वृत्तीय मांसपेशियां शिथिल होती हैं, रेडियल मांसपेशियां संकुचित होती हैं और पुतलियां फैलती हैं।**
C.  वृत्तीय और रेडियल मांसपेशियां संकुचित होती हैं जबकि पुतलियां फैलती हैं।
D.  वृत्तीय मांसपेशियां संकुचित होती हैं, रेडियल मांसपेशियां शिथिल होती हैं और पुतलियां संकुचित होती हैं।

**Row 2**

Language: Hindi

Question Image:



Question and Answer:

[JEE (Main)] [Physics]

**Question:** चित्र (a), (b), (c), (d) देखकर निर्धारित करें कि ये चित्र क्रमशः किन सेमीकंडक्टर डिवाइसों के अभिलक्षणांक ग्राफ हैं :
**Options:**
**A. साधारण डायोड, जीनर डायोड, सोलर सेल, LDR (लाइट डिपेंडेंट रेजिस्टेंस)**
B. जीनर डायोड, साधारण डायोड, LDR (लाइट डिपेंडेंट रेजिस्टेंस), सोलर सेल
C. सोलर सेल, LDR (लाइट डिपेंडेंट रेजिस्टेंस), जीनर डायोड, साधारण डायोड
D. जीनर डायोड, सोलर सेल, साधारण डायोड, LDR (लाइट डिपेंडेंट रेजिस्टेंस)

**Row 3**

Language: Hindi

Question Image:



Question and Answer:

[UP-CET] [Social Sciences]

**Question:** चित्र में दिए किले को पहचानिये।
**Options:**
A. जोधपुर किला
**B.  ग्वालियर किला**
C. लाल किला
D. आमेर किला

**Row 4**

Language: Hindi

Question Image:

$$\int_1^2 (x^2 - 2x + 4)^{3/2}\, dx = \frac{k}{k+5}$$

Question and Answer:

[JEE (Main)] [Mathematics]

**Question:** यदि इस छवि में दिखाए गए समीकरण के अनुसार, तो $k$ बराबर है :
**Options:**
A. 1
B. 2
**C. 3**
D. 4

| Language | Question Image | Question and Answer |
|---|---|---|
| Hindi |  | <br>**Question:** विभिन्न देशों की औद्योगिक वृद्धि (₹ करोड़ में) में निम्नलिखित में से कितने देशों की औद्योगिक वृद्धि औसत औद्योगिक वृद्धि से अधिक है?<br>**Options:**<br>A. 4<br>**B. 2**<br>C. 3<br>D. 1 |
| Lithuanian |  | High School Exam — Geography<br>**Question:** Kiek platumos laipsnių yra tarp Šiaurės poliarinio rato ir Pietų atogrąžos?<br>**Options:**<br>A. Apie 23°.<br>**B. Apie 90°.**<br>C. Apie 100°.<br>D. Apie 132°. |
| Nepali |  | PSC Exam — Reasoning<br>**Question:** दिइएको चित्र १,२,३,४ र ५ मध्येबाट कुनै तीन चित्रहरु एक आपसमा मिलाउदा पुर्ण आकारको त्रिभुजको चित्र बन्दछ । उक्त पूर्ण आकारको बनाउने चित्रहरुको नम्बरहरु दिइएको विकल्पबाट छनौट गर्नुहोस् ।<br>**Options:**<br>A. 124<br>**B. 234**<br>C. 245<br>D. 345 |
| Persian |  | Olympiad of Informatics — Math<br>**Question:**<br>طبق شکل زیر تعدادی چرخ‌دنده داریم که با هم درگیر هستند. چند دور و در کدام جهت باید چرخ دنده‌ی b را بچرخانیم تا چرخ‌دنده‌ی a دقیقا یک دور ساعت‌گرد بچرخد؟ تعداد دنده‌های چرخ‌دنده‌ی کوچک ۸، چرخ‌دنده‌های متوسط ۱۶ و چرخ‌دنده‌های بزرگ ۳۲ است.<br>**Options:**<br>A، 1 دور ساعت‌گرد<br>B، 1 دور پادساعت‌گرد<br>C، 2 دور ساعت‌گرد<br>**D، نمی‌توان چرخ‌دنده‌ی a را چرخاند** |

| Language | Question Image | Question and Answer |
|---|---|---|

**Row 1:**

Persian



Driving Test | Driving

**Question:**

عنوان این تابلو چیست؟

**Options:**

A، پیچ های پی در پی (اولین پیچ به چپ)

**B، پیچ های پی در پی (اولین پیچ به راست)**

C، پیچ به راست

D، پیچ به چپ

---

**Row 2:**

Persian



High School Exam | Engineering

**Question:**

شکل مقابل نمایانگر کدام است؟

**Options:**

A، قانون دوم کپلر

B، (نظریهٔ مه‌بانگ)

**C، کهکشان راه شیری**

D، راه مکه

---

**Row 3:**

Persian



University Entrance Exam | Physics

**Question:**

گلوله‌ای مسیری مطابق شکل را طی می‌کند. کار نیروی وزن از A تا B چند برابر کار نیروی وزن از B تا C است؟

**Options:**

A. 1.5
B. 1.25
C. 1.2
D. 2

---

**Row 4:**

Persian



University Entrance Exam | Geology

**Question:**

با توجه به گزینه‌ها، در شکل روبرو به ترتیب قدیمی‌ترین و جوان‌ترین پدیده کدام است؟

**Options:**

A، پیشروی دریا – گسل عادی

**B، پسروی دریا – تزریق دایک**

C، پیشروی دریا – ناپیوستگی هم‌شیب

D، پسروی – گسل عادی

| Language | Question Image | Question and Answer |
|---|---|---|
| Portuguese |  | **UNESP**  **Social Sciences**<br><br>**Question:** Observe as fachadas de duas igrejas. À esquerda, a Basílica de San Michele, construída no século XII em Pavia, na Itália. À direita, a Catedral de Reims, erguida a partir do século XIII em Reims, na França.<br>(Georges Duby e Michel Laclotte (orgs.). História artística da Europa: a Idade Média II, 1998.)<br>As duas fachadas<br>**Options:**<br>**A. diferenciam-se pela pouca ornamentação de San Michele, que expressa o estilo românico, e pela monumentalidade e sofisticação de Reims.**<br>B. diferenciam-se pela solidez de San Michele, que simboliza a força espiritual do catolicismo, e pela carência de detalhes na sede papal em Reims.<br>C. igualam-se na suntuosidade e no rebuscamento arquitetônico, indicando o poderio econômico da Igreja católica.<br>D. diferenciam-se pela discrição de San Michele, que revela o rigor na conduta dos protestantes, e pela ostentação da riqueza católica de Reims. |
| Portuguese |  | **FAMERP Entrance Exam**  **Physics**<br><br>**Question:** Quando um gerador de força eletromotriz 12 V é ligado a um resistor R de resistência $5,8\Omega$, uma corrente elétrica i de intensidade 2,0 A circula pelo circuito.<br>R<br>A resistência interna desse gerador é igual a<br>**Options:**<br>A. $0,40\Omega$.<br>**B. $0,20\Omega$**<br>C. $0,10\Omega$.<br>D. $0,30\Omega$. |
| Portuguese |  | **Unicamp Entrance Exam**  **Language**<br><br>**Question:** A imagem a seguir apresenta a transcrição de um diálogo em um vídeo publicado no Instagram.<br>No diálogo, a principal característica da reformulação da fala da médica é a inserção de<br>**Options:**<br>A. expressões que utilizam verbos frasais para recontextualizar o tratamento da paciente.<br>B. abreviações de substantivos, através das quais a médica amplia as informações do caso.<br>C. gírias que utilizam diversas classes de palavras para especificar melhor o diagnóstico da paciente.<br>**D. vocábulos marcados pela oralidade, através dos quais a médica atualiza os procedimentos futuros.** |

| Language | Question Image | Question and Answer |
|---|---|---|
| Portuguese |  | **ENEM, Brazil**  **Mathematics**<br><br>**Question:** Um segmento de reta está dividido em duas partes na proporção áurea quando o todo está para uma das partes na mesma razão em que essa parte está para a outra. Essa constante de proporcionalidade é comumente representada pela letra grega $\varphi$, e seu valor é dado pela solução positiva da equação $\varphi^2 = \varphi + 1$.<br>Assim como a potência $\varphi^2$, as potências superiores de $\varphi$ podem ser expressas da forma $a\varphi + b$, em que a e b são inteiros positivos, como apresentado no quadro.<br>A potência $\varphi^7$, escrita na forma $a\varphi + b$ ( a e b são inteiros positivos), é<br>**Options:**<br>A. $7\varphi + 2$<br>B. $9\varphi + 6$<br>C. $11\varphi + 7$<br>**D.** $13\varphi + 8$ |
| Serbian |  | **Mathematical Kangaroo**  **Mathematics**<br><br>**Question:** Колико процената површине троугла на слици је осенчено?<br>**Options:**<br>**A. 88 %**<br>B. 90 %<br>C. 85 %<br>D. 80 % |
| Russian |  | **Mathematical Kangaroo**  **Mathematics**<br><br>**Question:** Каких геометрических фигур нет на рисунке?<br>**Options:**<br>A. кругов<br>B. все эти фигуры есть<br>C. прямоугольников<br>**D. треугольников** |
| Spanish |  | **Medicine Exam**  **Pulmonology**<br><br>**Question:** Varón de 60 años, fumador activo, que presenta tos y expectoración diaria de años de evolución, ocasionalmente hemoptoica. En los últimos meses se añade disnea progresiva. Presenta acropaquia y en la auscultación pulmonar destacan roncus y sibilantes teleinspiratorios en pulmón izquierdo. La TC pulmonar de alta resolución se muestra en la imagen adjunta. ¿Cuál es el diagnóstico más probable?<br><br>**Choices:**<br>A. Carcinoma quístico.<br>B. Enfisema pulmonar.<br>C. Tuberculosis cavitada.<br>**D. Bronquiectasias.** |

The table in the Portuguese image:

| $\varphi^2$ | $\varphi^3$ | $\varphi^4$ | $\varphi^5$ | $\varphi^6$ | $\varphi^7$ |
|---|---|---|---|---|---|
| $\varphi + 1$ | $2\varphi + 1$ | $3\varphi + 2$ | $5\varphi + 3$ | $8\varphi + 5$ | ... |

| Language | Question Image | Question and Answer |
|---|---|---|
| Spanish |  | **Undergraduate Exam** · **Biophysics**<br><br>**Question:** Calcule el valor de la primera resistencia (R1)<br>**Options:**<br>A. 42 $\Omega$<br>B. 6 $\Omega$<br>**C. 12** $\Omega$<br>D. 24 $\Omega$ |
| Spanish |  | **High School Exam, Colombia** · **Biology**<br><br>**Question:** En un laboratorio se estudia el comportamiento del volumen de un gas ideal al variar su temperatura, obteniendo la siguiente gráfica: Teniendo en cuenta la información de la gráfica, si la temperatura aumenta de -153 °C a -33 °C, ¿qué pasa con el volumen del gas?<br>**Options:**<br>A. Disminuye de 30 L a 25 L.<br>B. Disminuye de 10 L a 5 L.<br>C. Aumenta de 0 L a 10 L.<br>**D. Aumenta de 10 L a 20 L.** |
| Telugu |  | **Undergraduate Exam** · **Chemistry**<br><br>**Question:** ఇచ్చిన చిత్రంలో సమ్మేళనం యొక్క మోలార్ ద్రవ్యరాశి ఎంత?<br>**Options:**<br>A. 304.9<br>B. 304.4<br>C. 301.9<br>**D. 303.4** |
| Ukrainian |  | **ZNO Vision** · **Mathematics**<br><br>**Question:** Пластикові кульки радіуса 6 см зберігають у висувній шухлядці, що має форму прямокутного паралелепіпеда (див. рисунок). Якою з наведених може бути висота $h$ цієї шухлядки?<br>**Options:**<br>A. 3 см<br>B. 6 см<br>C. 10 см<br>**D. 13 см** |
| Ukrainian |  | **Driving Test** · **Driving**<br><br>**Question:** По якій траєкторії можна продовжити рух праворуч легковому автомобілю?<br>**Options:**<br>A. Тільки по А.<br>**B. Тільки по Б.**<br>C. По А і Б.<br>D. По будь-якій. |

Table 16: Samples from various exams in the KALEIDOSCOPE benchmark. The correct answer is highlighted in **Bold Green**. Some samples are reformatted for better presentation.