

# A Culturally-diverse Multilingual Multimodal Video Benchmark & Model

Bhuiyan Sanjid Shafique<sup>1\*</sup>, Ashmal Vayani<sup>2\*</sup>, Muhammad Maaz<sup>1</sup>, Hanoona Abdul Rasheed<sup>1</sup>, Dinura Dissanayake<sup>1</sup>, Mohammed Irfan Kurpath<sup>1</sup>, Yahya Hmaiti<sup>2</sup>, Go Inoue<sup>1</sup>, Jean Lahoud<sup>1</sup>, Md. Safirur Rashid<sup>3</sup>, Shadid Intisar Quasem<sup>3</sup>, Maheen Fatima<sup>4</sup>, Franco Vidal<sup>2</sup>, Mykola Maslych<sup>2</sup>, Ketan Pravin More<sup>1</sup>, Sanoojan Baliah<sup>1</sup>, Hasindri Watawana<sup>1</sup>, Yuhao Li<sup>1</sup>, Fabian Farestam<sup>5</sup>, Leon Schaller<sup>6</sup>, Roman Tymtsiv<sup>7</sup>, Simon Weber<sup>6</sup>, Hisham Cholakkal<sup>1</sup>, Ivan Laptev<sup>1</sup>, Shin'ichi Satoh<sup>8</sup>, Michael Felsberg<sup>10</sup>, Mubarak Shah<sup>2</sup>, Salman Khan<sup>1,9</sup>, Fahad Shahbaz Khan<sup>1,10</sup>

<sup>1</sup>Mohamed bin Zayed University of Artificial Intelligence, <sup>2</sup>University of Central Florida, <sup>3</sup>Islamic University of Technology, <sup>4</sup>Air University, <sup>5</sup>ETH Zurich, <sup>6</sup>Technische Universität München, <sup>7</sup>Independent Researcher, <sup>8</sup>National Institute of Informatics, <sup>9</sup>Australian National University, <sup>10</sup>Linköping University

{bhuiyan.shafique, muhammad.maaz, hanoona.bangalath, dinura.dissanayake, mohammedirfan.k, go.inoue, jean.lahoud, ketan.more, sanoojan.baliah, hasindri.watawana, hisham.cholakkal, salman.khan, fahad.khan}@mbzuai.ac.ae  
{ashmal.vayani, yohan.hmaiti, fr543419, mykola.maslych}@ucf.edu, {safirurrashid, shadidintisar}@iut-dhaka.edu  
{leon.schaller, simon.l.weber}@tum.de, {tymtsiv.roman, lyh88524}@gmail.com, maheen.fatima@students.au.edu.pk  
ffarestam@student.ethz.ch, satoh@nii.ac.jp, michael.felsberg@liu.se, ivan.laptev@inria.fr, mubarak.shah@crcv.ucf.edu

**Sinhala Examples:**

- Question:** මෙම සිංහල අලුත් අවුරුදු සම්ප්‍රදායේ අරමුණ කුමක්ද?
- Answers:**
  - Vi:** සෑම පැත්තකින්ම කිරි නිසි ලෙස ගලා එන්නේ නම්, එය සමෘද්ධිමත් නම් වසරක් සඳහා හොඳ සලකුණක් ලෙස සැලකේ.
  - ELM:** The video shows a man cooking food in a pot over an open fire. The pot is filled with rice and milk, and the man is stirring the contents with a spoon. The video is in Tamil.
  - EV:** The Sinhala New Year is celebrated by cooking rice in a traditional clay pot over a fire made from wood.
  - OQ:** මෙම සිංහල අලුත් අවුරුදු සම්ප්‍රදායේ අරමුණ කුමක්ද?
  - LN:** සුදු අලුත් අවුරුදු සම්ප්‍රදායේ අරමුණ කුම දී ගැනයාට අපේ අතුරේ අතුරේ.
  - ELM:** සම්ප්‍රදායේ අරමුණ කුමක්ද අවුරුදු සම්ප්‍රදායේ අරමුණ කුමක්ද.

**Bengali Examples:**

- Question:** ভিডিও থেকে গরুর মাংসের তেহারিটির স্বাদ বর্ণনা করুন।
- Answers:**
  - Vi:** ভিডিওতে গরুর মাংসের তেহারিটি স্বাদময় এবং ভালো মনে হচ্ছে।
  - ELM:** ভিডিও থেকে গরুর মাংসের তেহারিটির স্বাদ বর্ণনা করুন।
  - EV:** The video does not provide any information about the taste of the garur mangsho tehari.
  - OQ:** ভিডিও থেকে গরুর মাংসের তেহারিটির স্বাদ বর্ণনা করুন।
  - LN:** ভিডিও থেকে গরুর মাংসের তেহারিটির স্বাদ বর্ণনা করা যাচ্ছে। এর কারণ ভিডিও থেকে গরুর মাংসের তেহারিটির স্বাদ বর্ণনা করা যাচ্ছে।
  - ELM:** ভিডিও থেকে গরুর মাংসের তেহারিটির স্বাদ বর্ণনা করুন।

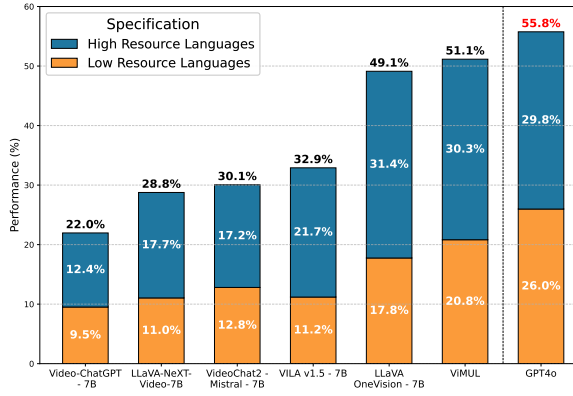
Figure 1: ViMUL-Bench consists of carefully curated videos spanning 14 languages, with 8K manually verified annotations by native experts. It covers 15 diverse domains, incorporating real-world cultural elements such as regional landmarks, local cuisines, and traditional festivals. Additionally, we introduce ViMUL, a simple multilingual baseline designed for general and cultural video comprehension. Qualitative examples (top: Sinhala and bottom: Bengali language) here show that ViMUL performs favorably against recent vidLMs in cultural inclusivity and overall understanding (errors are highlighted in red and correct answer in green). ViMUL-Bench covers diverse questions, such as MCQs and short and long visual question answers (VQAs). (ELM: ViLA, EV: Video-Chat2, Video-ChatGPT, OQ: LLaVA-OneVision-Qwen (OQ), LN: LLaVA-Next (LN), Vi: Our ViMUL).

## Abstract

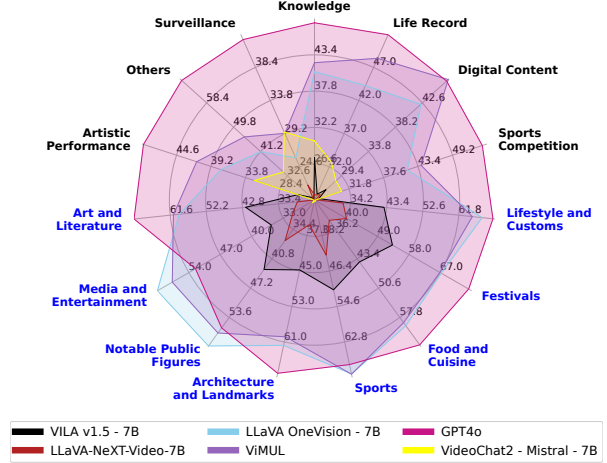
Large multimodal models (LMMs) have recently gained attention due to their effectiveness to understand and generate descriptions of visual content. Most existing LMMs are in English language. While few recent works explore multilingual image LMMs, to the best of our knowledge, moving beyond the En-

glish language for cultural and linguistic inclusivity is yet to be investigated in the context of video LMMs. In pursuit of more inclusive video LMMs, we introduce a multilingual Video LMM benchmark, named ViMUL-Bench, to evaluate Video LMMs across 14 languages, including both low- and high-resource languages: Arabic, Bengali, Chinese, English, French, German, Hindi, Japanese, Russian, Sinhala, Spanish, Swedish, Tamil, and Urdu.

\*Equal contribution



(a) Open vs. Closed-Source LMM performance.



(b) Performance on various general and cultural categories.

Figure 2: **Benchmarking video LMMs on the proposed ViMUL-Bench across various languages and cultures.** (a) Performance comparison of open-source versus closed-source models, with a distinction between low-resource and high-resource languages in our ViMUL-Bench. (b) Performance of different video LMMs across 15 diverse categories (both generic and cultural) in our ViMUL-Bench. The categories in **black** represents generic categories, and categories in **blue** represents the cultural categories.

Our ViMUL-Bench is designed to rigorously test video LMMs across 15 categories including eight culturally diverse categories, ranging from lifestyles and festivals to foods and rituals and from local landmarks to prominent cultural personalities. ViMUL-Bench comprises both open-ended (short and long-form) and multiple-choice questions spanning various video durations (short, medium, and long) with 8k samples that are manually verified by native language speakers. In addition, we also introduce a machine translated multilingual video training set comprising 1.2 million samples and develop a simple multilingual video LMM, named ViMUL, that is shown to provide a better tradeoff between high-and low-resource languages for video understanding. We hope our ViMUL-Bench and multilingual video LMM along with a large-scale multilingual video training set will help ease future research in developing cultural and linguistic inclusive multilingual video LMMs. Our proposed benchmark, video LMM code-base and training data are available at <https://mbzuai-oryx.github.io/ViMUL/>.

## 1 Introduction

Large Multimodal Models (LMMs) have achieved remarkable success in vision-and-language tasks, yet their development has predominantly centered on English, overlooking the vast linguistic and cultural diversity of global users (Vayani et al., 2025; Pfeiffer et al., 2021). This English-centric focus leads to significant performance gaps for other languages, as models often fail to grasp user intent

when queries or captions are in low-resource languages. Moreover, current models struggle with cultural nuances and region-specific context specifically for low-resource languages (Romero et al., 2024; Qureshi et al., 2025; Raza et al., 2025a). For instance, the MaRVL (Liu et al.) and the recent ALM-Bench (Vayani et al., 2025) image LMM benchmarks, verified by native speakers to include diverse (including low-resource) languages, reveal dramatic drops in accuracy when state-of-the-art models operate beyond English. These findings underscore the pressing need for multilingual and multicultural evaluation benchmarks to develop more inclusive next generation of LMMs.

Existing efforts to explore linguistically and culturally diverse LMM benchmarks are limited to *images* (Romero et al., 2024; Vayani et al., 2025). To the best of our knowledge, linguistic and cultural diversity are yet to be investigated for *video* LMMs. Video domain poses different challenges as it often depicts complex, culturally rich scenarios-local festivals, foods, rituals, or landmarks, that require understanding both the visual context and the language-specific narration or questions. Fig. 1 illustrates an example where LMMs are asked to describe the taste of the Bengali dish *Beef Tehari*. The models fail to interpret the question in the local language, responding incorrectly and missing linguistic nuances. While short and long video understanding LMM benchmarks exist in literature, they are typically restricted to only English language. For instance, Video-MME (Fu et al., 2024) focuses

on diverse video analysis but in a single language, and MVBench (Li et al., 2024b) emphasizes temporal reasoning (action sequences, motion) without multilingual considerations. Other recent efforts like ViLMA (Kesen et al., 2024) and SEED-Bench (Li et al., 2023b,a) probe video-language model’s abilities in zero-shot temporal grounding and procedural understanding, among other skills, yet none assess cross-lingual or cross-cultural comprehension. In short, there is no comprehensive benchmark to evaluate how well video LMMs perform across different languages and cultural contexts (see Tab. 1).

To bridge this gap, we propose Multilingual Video LMM Benchmark (ViMUL-Bench), the first benchmark for evaluating video LLMs across 14 languages spanning both high-resource and low-resource cases. Besides being multilingual, our ViMUL-Bench is designed to test cultural awareness in video LMMs. It covers a broad spectrum of culturally diverse categories, including distinct lifestyles, traditional festivals, local cuisines, rituals, regional landmarks, and notable cultural figures. We formulate a rich evaluation suite with both open-ended questions (requiring descriptive answers in short or long form) and multiple-choice questions (MCQs), curated for videos of varying lengths (short clips, medium snippets, and longer videos) to assess understanding at different temporal scales. Crucially, the entire benchmark is verified by native speakers of each language, ensuring that questions and answers accurately capture nuances of tradition, customs, and societal context. Additionally, we construct a specialized multilingual video training dataset and train a strong baseline model named ViMUL on it. Our experimental analysis reveals that ViMUL provides a better tradeoff between high- and low-resource languages, achieving superior overall performance on multilingual multicultural video question answering, compared to existing open-source video LMMs. Our contributions are summarized as follows:

- We introduce ViMUL-Bench, a comprehensive benchmark for video LMMs covering 14 languages (including several under-represented ones) and 15 domains, including real-world cultural aspects. To our knowledge, this is the first effort to enable rigorous testing of video LMMs across a wide linguistic and cultural spectrum, emphasizing both cross-lingual and cultural comprehension (see

Tab. 1).

- The ViMUL-Bench offers 8K manually verified diverse samples for comprehensive spatio-temporal evaluation and includes both open-ended and multiple-choice QAs. It also offers diversity in terms of video length (short, medium, and long videos). In addition, we provide a multilingual training dataset with 1.2M samples translated from available video datasets.
- We propose ViMUL, a multilingual video LLM fine-tuned on our multilingual training set. ViMUL establishes a strong baseline on ViMUL-Bench, providing a better overall tradeoff compared to existing open-source video LMMs for multilingual video understanding (see Fig. 2).

## 2 Related Works

**Multilingual Multicultural Datasets:** Early vision-language benchmarks were predominantly English-centric, with limited coverage of other languages or cultures (Romero et al., 2024). Recent efforts have sought to bridge this gap by extending multimodal tasks to multiple languages. For example, xGQA expanded the GQA visual question answering dataset to seven diverse languages via translation (Pfeiffer et al., 2021). However, such translation-based approaches often reuse the same generic/biased image distributions and thus fail to capture cultural nuances (Romero et al., 2024). To introduce culture-specific content, (Liu et al.) proposed a Multicultural Reasoning dataset (MaRVL). While MaRVL incorporates diverse concepts, its scope is limited (five languages and binary true/false reasoning). Similarly, other contemporary benchmarks include M3Exam (Zhang et al., 2023b), MMBB (Sun et al., 2024), MM-Bench (Liu et al., 2024), M4U (Wang et al., 2024), and Exams-V (Das et al., 2024) which introduce multilingual evaluation samples for image understanding. However, their scope is generally limited to a few languages and offers a narrow cultural scope. More recent benchmarks extend the scope further, e.g., CVQA (Romero et al., 2024) introduces 10K VQA examples grounded in 31 languages (13 scripts), using culturally relevant images and human-curated questions. The recent All Languages Matter (ALM-Bench) benchmark (Vayani et al., 2025) spans 100 languages with im-

Benchmark	Multilingual	Total Domains	Total Samples	Question Types	Total Videos	Annotation Type	Cultural Content
ActivityNet-QA							
(Fabian Caba Heilbron and Niebles, 2015)	✗ (1)	-	2378	OE	5800	Human	✗
CVRR-ES (Khattak et al., 2024)	✗ (1)	-	2400	OE	224	Auto	✗
MoVQA (Zhang et al., 2023a)	✗ (1)	6	21,953	OE	-	Human	✗
MovieQA (Tapaswi et al., 2016)	✗ (1)	-	6462	MCQ	6771	Human	✗
MSVD-QA (Xu et al., 2017)	✗ (1)	5	50,505	OE	1970	Auto	✗
MVBench (Li et al., 2024b)	✗ (1)	-	4000	MCQ	3507	Auto	✗
Perception Test (Patraucean et al., 2023)	✗ (1)	-	44,000	MCQ+OE	11,600	Auto+Human	✗
TVQA (Lei et al., 2018)	✗ (1)	-	15,2545	MCQ+OE	21,793	Human	✗
Video-MME (Liu et al.)	✗ (1)	6	2700	MCQ	900	Human	✗
VCG Diverse (Maaz et al., 2024b)	✗ (1)	18	4354	SVQA, LVQA	877	Auto	✗
<b>Ours</b>	✓ (14)	15	8,025	MCQ, SVQA, LVQA	879	Auto+Human	✓

Table 1: Comparison of video LMM benchmarks emphasizing multilingual and cultural understanding. *Domains* represent the aspects covered by each dataset for different languages. *Annotation Type* is categorized as follows: Human - Questions were created in the local language. Human+Auto - Questions were generated or translated using GPT-4/Google API and later validated by human experts. Auto: Questions were generated or translated automatically without human validation. ‘-’ indicates that information is not available.

ages drawn from 13 distinct cultural aspects. It is the largest multicultural multimodal evaluation benchmark, designed to test LMMs on diverse imagery and low-resource languages. However, these benchmarks remain limited to the image-domain and do not address multilingual and multicultural aspects unique to videos.

**Video LMM Benchmarks:** Extending multimodal evaluation to video introduces additional challenges due to temporal dynamics. A number of video-language benchmarks have emerged to assess LMMs on video understanding, though they focus on general monolingual capabilities. Video-MME (Fu et al., 2024) is introduced as the first comprehensive evaluation of LMMs on video analysis, covering a broad spectrum of video domains. It offers a full-spectrum evaluation, covering questions from perception to reasoning, and variable videos lengths ( $\sim 10$  secs to  $\sim 1$  hour). Another effort, MVBench (Li et al., 2024b) concentrates on temporal reasoning skills and proposes an automatic pipeline to generate a large multiple-choice QA benchmark by leveraging existing video datasets and GPT-based annotators. Beyond general benchmarks, ViLMA (Kesen et al., 2024) takes a fine-grained approach as it uses carefully designed counterfactual video pairs to probe a model’s temporal grounding and linguistic understanding in a zero-shot setting. SEED-Bench is another comprehensive multimodal benchmark that includes some video-based questions; however, those are largely confined to temporal or procedural understanding tasks (Li et al., 2023b,a). These benchmarks have a significantly advanced evalua-

tion for video-based multimodal reasoning and generation, however they lack multilingual and culturally aware components. Current video benchmarks evaluate models primarily on English video-question pairs and generic content, without testing performance on non-English dialogues, region-specific contexts, or culturally diverse visual narratives. The absence of multilingual and multicultural evaluation for video LMMs forms a key motivation for ViMUL-Bench, which aims to fill that gap by assessing video understanding across diverse languages and cultural settings.

### 3 ViMUL-Bench

ViMUL-Bench is a comprehensive multilingual benchmark, designed to evaluate both general and culturally-specific aspects of video comprehension in video LMMs. It captures rich cultural nuances through a diverse set of question-answer (QA) pairs, including multiple-choice and open-ended (short and long) formats (Raza et al., 2025b; Campos et al., 2025). The benchmark spans 15 diverse categories, categorized into general and cultural topics, across 14 languages: Arabic, Bengali, Chinese, English, French, German, Hindi, Japanese, Russian, Sinhala, Spanish, Swedish, Tamil, and Urdu. We build upon the recent PALO (Maaz et al., 2025), incorporating its 10 languages while adding Swedish, German, Tamil, and Sinhala to ensure typological diversity and to enhance the representation of low-resource languages such as Tamil, Urdu, and Sinhala as defined by Ethnologue (Campbell and Grondona, 2008) and Glottolog (Hammarström et al., 2022) database.

The *generic* category includes seven domains: Artistic Performance, Digital Content, Knowledge, Life Record, Sports Competition, Surveillance, and Others. The *cultural* categorization is inspired from recent image LMM benchmarks (Vayani et al., 2025; Romero et al., 2024; Marino et al., 2019), where the corresponding videos for each domain and language are manually scrapped with their annotations manually verified by a native speaker. It spans eight diverse domains, including Lifestyle & Customs, Festivals, Food & Cuisine, Sports, Architecture & Landmarks, Notable Public Figures, Media & Entertainment, and Art & Literature. Our ViMUL-Bench has been meticulously curated and verified by native experts for the 13 languages to ensure high-quality question answer (QA) pairs that accurately capture the nuances of all 15 domains. It

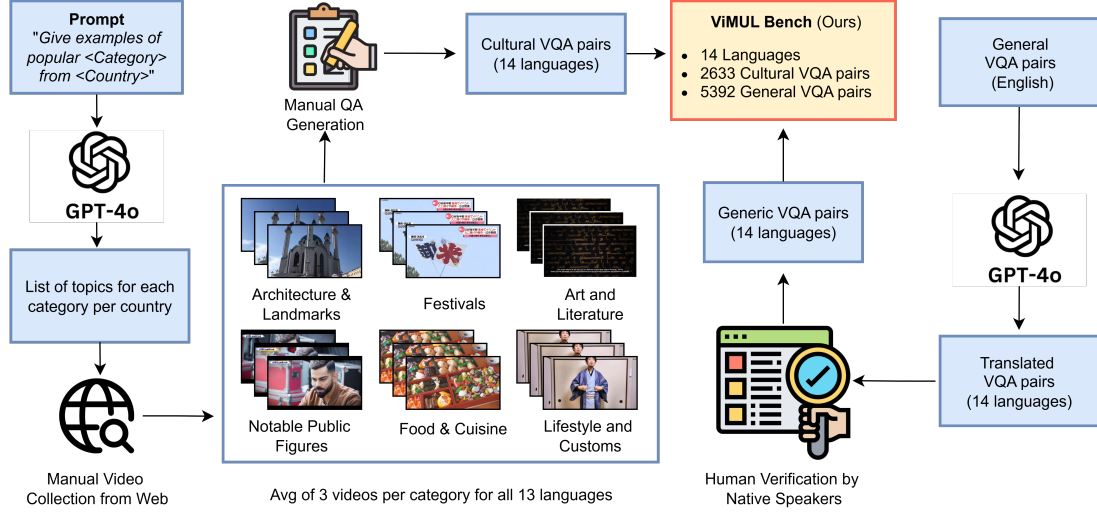


Figure 3: **Data collection and verification pipeline.** Our benchmark consists of both cultural-specific video content curated from scratch (*left*) and generic Video-QA pairs sourced from existing video LMM benchmarks. Cultural videos are scrapped using a (*country, language, sub-topic*) triplet and manually filtered for relevance and private information. With the help of native speakers, we create QA pairs for each language from scratch (except English), with cultural QA pairs translated into English using GPT-4o. Our ViMUL-Bench has diverse question types and features approximately 8K QA pairs in 14 languages.

comprises 8,025 diverse questions in total across 14 languages, spanning both generic and cultural categories to comprehensively evaluate multilingual and cross-cultural video understanding.

### 3.1 Data Collection and Annotation

**Generic VQA Curation.** As discussed above, our ViMUL-Bench encompasses both generic and cultural categories to evaluate multilingual video understanding across 14 languages. For the generic category, we carefully curate English-only subsets from VCG-Diverse (Maaz et al., 2024b), CVRR-ES (Khattak et al., 2024), MVBench (Li et al., 2024b), and VideoMME (Fu et al., 2024). Among these, VCG-Diverse and CVRR-ES follow an open-ended question format, while MVBench and VideoMME use multiple-choice questions (MCQs). These English samples are then translated into 13 additional languages using GPT-4o, followed by manual verification and refinement by native speakers of each language. Language experts were explicitly instructed to correct poor translations, ensuring accuracy by rephrasing or rewriting question-answer pairs when necessary. This process resulted in 5,392 QA pairs from 542 videos, in the generic category (see Fig. 3).

**Cultural Video Curation.** To curate diverse cultural QA pairs, we collect open-licensed videos and their corresponding metadata from the internet, focusing on specific cultural aspects of each language across three durations: short (0-4 mins),

medium (4-15 mins), and long (15+ mins). To accurately capture the cultural nuances of each language, we follow the approach outlined in (Vayani et al., 2025; Romero et al., 2024) and link each language to a country based on the World Values Survey (Haerpfer et al., 2022) to ensure coverage of cultural and ritual diversity. Additionally, we generate a list of topics for each category, forming a triplet (language-country-topic), using GPT-4o, and then search for relevant content. For example, querying “Popular sports in United States” yields responses such as “baseball, soccer, golf, and ice hockey.” We manually extract videos in the native language to ensure linguistic accuracy. Each domain undergoes several filtration steps, such as removing low-resolution, noisy, or unclear videos, to guarantee data quality. To maintain both high-quality and cultural relevance, we enlist expert native speakers of each language to manually verify the quality and cultural diversity of videos. Any content lacking cultural relevance is removed from the dataset. Fig. 3 shows our data collection and verification pipeline.

**Cultural QA Generation.** To generate high-quality video-QA pairs for the cultural section of ViMUL-Bench, the question-answer (QA) pairs are curated via native experts based on the provided videos and their metadata. Notably, videos and their corresponding QAs are not shared across languages in the cultural set. For each video, we generate multiple-choice questions (MCQs) and

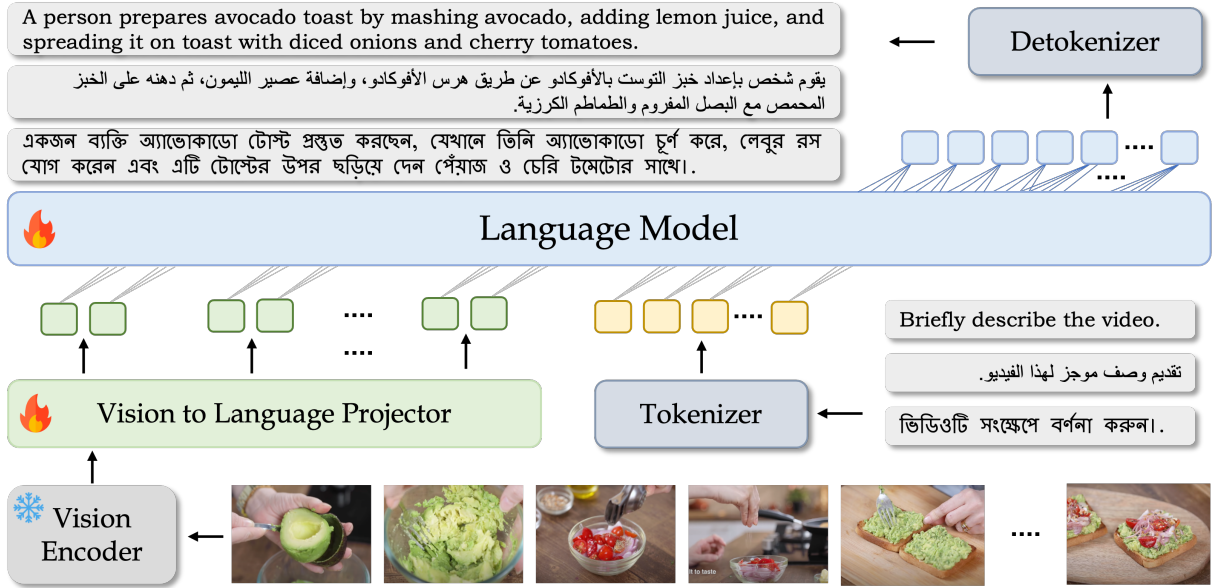


Figure 4: **Overview of ViMUL.** ViMUL is designed to comprehend and generate content in 14 different languages: Arabic, Bengali, Chinese, English, French, German, Hindi, Japanese, Russian, Sinhala, Spanish, Swedish, Tamil, and Urdu, covering at least two-thirds of the global population. The model employs a vision encoder to process video frames, followed by a vision-to-language projector and an LLM. The projected features are then concatenated with the user query and fed into the LLM to generate a response. (❄️: frozen, 🔥: trained)

one open-ended question in the native language. To reduce randomness in MCQ generation, we ensure that each question can also be answered when rephrased as an open-ended question. Further, the native experts are tasked with creating an English version of the question-answer pair. We instruct the experts to focus on cultural concepts depicted in each video and to generate questions that require a visual understanding of the video, while avoiding the perpetuation of bias and stereotypes. Following this process, 2,633 QA pairs were curated, spanning 337 videos for the cultural category.

#### 4 ViMUL: Multilingual Video LMM

In addition to the ViMUL-Bench, we develop a multilingual video LMM, named ViMUL, by constructing a multilingual video training set. ViMUL is built to understand and generate content in 14 diverse languages, covering an audience that represents at least two-thirds of the world’s population.

**Overall Architecture:** The architecture of ViMUL is derived from LLaVA-OneVision (Li et al., 2024a), which seamlessly integrates a vision encoder, a vision-to-language projector, and a language model. The video frames are encoded using a vision encoder, projected into the language model’s embedding space using a two-layer MLP projector, concatenated with the text embeddings, and passed to the language model to generate the

	Arabic	Bengali	Chinese	English	French	German	Hindi	Japanese	Russian	Sinhala	Spanish	Swedish	Tamil	Urdu
GPT-4o	55.0	53.2	55.7	63.2	58.3	54.9	51.6	50.4	58.7	49.6	55.3	55.9	50.1	49.6
ViMUL	50.6	47.5	49.0	61.8	58.4	53.0	48.4	48.0	55.8	31.5	54.6	50.0	32.7	47.0
LLaVA-OneVision-7B	51.0	41.0	53.2	65.1	54.1	51.4	41.3	46.3	53.6	24.0	54.4	45.0	28.2	37.8
VILA v1.5-7B	27.2	20.9	26.3	47.8	37.0	37.3	25.8	27.5	42.0	17.2	42.0	33.5	17.6	20.8
VideoChat2-Mistral-7B	30.1	26.8	28.8	31.4	32.2	34.8	31.9	29.8	36.0	20.4	32.0	33.2	21.3	28.4
LLaVA-NeXT-Video-7B	22.4	21.1	28.3	39.4	27.7	32.5	23.0	26.7	33.0	18.5	35.0	29.0	16.8	22.3
Video ChatGPT-7B	18.2	21.1	23.4	22.6	21.7	23.2	20.3	24.4	25.6	15.1	26.6	25.8	17.7	19.7

Figure 5: **Performance comparison of video LMMs across 14 languages on ViMUL-Bench.** Average accuracy is reported across all question types for each language. Each box represents a model’s accuracy for a specific language, with darker shades indicating higher accuracy. The results show that the closed-source model, GPT-4o, generally outperforms its open-source counterparts. In contrast to high-resource languages, methods struggle on low-resource languages (e.g., Sinhala, Urdu, Tamil). Among open-source models, our ViMUL provides a better tradeoff between high and low-resource languages, achieving an overall gain of 2% over LLaVA-OneVision.

response (see Fig. 4). We discuss the Video Sampling strategy in Sec. J (suppl. material).

#### 4.1 Multilingual Instruction Tuning Dataset

One of the contributions of this work is the development of a comprehensive multilingual video-language instruction-tuning dataset. Current video instruction-tuning datasets exist only in English and do not focus on other languages. However, recent advances in large language models (LLMs) have demonstrated impressive per-

formance in multilingual tasks. We leverage these advancements and use GPT-4o-mini (OpenAI, 2024) to translate the video instruction-tuning dataset from English to 13 additional languages, thereby creating a multilingual dataset that broadens the linguistic scope and applicability of the model. Our dataset is sourced from two primary sources: Video-Instruct100K (Maaz et al., 2024a) and LLaVA-Video-178K (Zhang et al., 2024). Video-Instruct100K is a video instruction-tuning dataset containing 100K samples generated using a semi-automatic annotation pipeline. The QA pairs are open-ended, including both short and long question-answer formats. We use the human-verified version of Video-Instruct100K released by VideoGPT+ (Maaz et al., 2024b), which consists of 25,803 samples. The videos in Video-Instruct100K are sourced from the ActivityNet dataset (Fabian Caba Heilbron and Niebles, 2015). Finally, our pipeline results in a total of 1,238,102 samples across all languages. We also discuss dataset-wise, per-language QA-pair distribution, and common translation issues in Sec. J (suppl. material).

To evaluate the translation quality generated by GPT-4o, we perform a *cycle consistency check* (Hu et al., 2011). We randomly sample 10,000 QA pairs across 10 languages, translate them back to English using Qwen-3 (Yang et al., 2025), and compare the results with our original English subset using GPT-4o as the judge. Fig. 16 (suppl. material) shows the per-language translation scores. The performance ranges from 95.3% in French to 84.4% in Bengali, demonstrating the quality of our multilingual data.

## 5 Results and Discussions

As discussed earlier, ViMUL-Bench comprises two question types: multiple-choice and open-ended, including both short and long question-answer formats, different prompts are employed for each question type. For multiple-choice questions, we provide the visual context and textual query to the LMMs, instructing them to select the best option, which is then directly compared to the ground truth. Performance is measured using *accuracy*, following established multiple-choice VQA benchmarks (Romero et al., 2024; Bang et al., 2023; Zhu et al., 2016). For open-ended questions, we use the open-source multilingual LLM, Phi-4-14B (Abdin et al., 2024) as a judge, ensuring consistency and reproducibility, unlike GPT-based models (Shen et al., 2023; Stureborg et al., 2024), which are costly and

inconsistent due to version updates. Performance is evaluated using *correctness* criteria, which measure how closely the model’s output matches the ground truth (see Sec. A in suppl. material for further detail). This approach follows recent work in evaluation frameworks (Vayani et al., 2025; Narnaware et al., 2025), though the *correctness* used here is specific to our setup.

Fig. 5 shows the per-language performance comparison of different video LMMs on ViMUL-Bench. The closed-source proprietary model, GPT-4o (OpenAI, 2024), consistently outperforms open-source models. Among open-source models, our multilingual VidLMM, ViMUL, achieves a better tradeoff with respect to high-and low-resource languages with an overall accuracy of 51.1%, followed by LLaVA-OneVision (Li et al., 2024a) with 49.1%. Both open-source and closed-source models face challenges with several low-resource languages, such as Sinhala, Tamil, and Urdu. For example, GPT-4o’s performance drops significantly from 63.2% on English to 49.6% on Urdu. Similarly, LLaVA-OneVision (Li et al., 2024a) drops from 65.1% on English to 24% on Sinhala, indicating that the model struggles with under-represented languages. In comparison, ViMUL outperforms LLaVA-OneVision by approximately 9.2% on Urdu, and 7.5% on Sinhala, showing its efficacy on these low-resource languages.

Fig. 2a presents the performance breakdown of video LMMs on low-resource and high-resource languages. ViMUL-Bench includes three low-resource languages (Sinhala, Urdu, and Tamil), as defined in Costa-Jussà et al. (2022). The results show that the performance gap between open-source and closed-source GPT-4o increases on low-resource languages.

**Effect of Question Type.** As mentioned earlier, ViMUL-Bench includes two types of questions: multiple-choice (MCQs) and open-ended (OE), with the latter further divided into long and short VQAs. Fig. 7 shows the performance of video LMMs on these question types. Overall, video LMMs perform better in MCQs but struggle to generate *correct* responses for OE questions. This is likely because OE questions are more complex and require enhanced understanding and multilingual reasoning across both generic and cultural domains. The only exception is VideoChat2 (Li et al., 2023c) and VideoChatGPT (Maaz et al., 2024a), which perform better on OE questions than on MCQs. Among all models, the closed-source GPT-




Sinhala: Sinhalese Script	Urdu: Arabic Script	Bengali: Bengali Script
 <p><b>Question:</b> සමරනු ලබන්නේ කුමක් ද?  <b>Ground Truth Answer:</b> බුදුරජාණන් වහන්සේගේ උපත, බුද්ධත්වය සහ පරිනිර්වාණය  <b>Predicted Answer:</b> සමරනු ලබන්නේ වෙසක් උත්සවය හෝ බුදුන්ගේ ජාතික, සම්බුද්ධත්ව සහ පරිනිර්වාණ දිනයයි.  <b>Error Type:</b> Language Error  <b>Video Duration:</b> Short (&lt;1 Min)</p>	 <p><b>Question:</b> ویشو کے آغاز میں کون سی مشہور ڈش دکھائی گئی ہے؟  (A. مچھر, B. باری کيو پليٹر, C. کھنہ, D. بریانی)  <b>Ground Truth Answer:</b> باری کيو پليٹر  <b>Predicted Answer:</b> C  <b>Error Type:</b> Lack of Cultural Understanding  <b>Video Duration:</b> Short (1+ Min)</p>	 <p><b>Question:</b> বাংলাদেশ কীভাবে গোল করল ?  <b>Ground Truth Answer:</b> শট গুই ইন থেকে বেড গোলটি করে বাংলাদেশ  <b>Predicted Answer:</b> বাংলাদেশ কর্নার কিং থেকে বেডের মাধ্যমে গোল করছে।  <b>Error Type:</b> Lack of Knowledge  <b>Video Duration:</b> Short (2+ Mins)</p>
Category: Festivals	Category: Food & Cuisine	Category: Sports

Figure 6: We present qualitative examples of failure cases of GPT-4o’s across different language scripts and categories, specifying the corresponding error types. For instance, in a Sinhala-language question asking about the event being celebrated in the video, the model correctly identifies the cultural significance-celebrating the birth, enlightenment, and passing of Lord Buddha, but fails to respond with grammatically correct Sinhala, highlighting a language proficiency error. Results on success cases are shown in Fig. 9 (suppl. material).

4o achieves the highest accuracy on OE questions (54.6%), followed by ViMUL, which leads among open-source models with 36%. Notably, ViMUL achieves the highest accuracy on MCQs (62.8%) among all methods. We also demonstrate the consistency of Phi-4 scores with human judgement on GPT-4o outputs in Sec. F (suppl. material).

**Performance across Language Scripts.** We group the 14 languages in our ViMUL-Bench by language scripts (Costa-Jussà et al., 2022), using data from Ethnologue (Campbell and Grondona, 2008) and the Glottolog database (Hammarström et al., 2022). This results in nine distinct scripts. Fig. 8 shows the performance of video LMMs across these scripts. The closed-source GPT-4o consistently achieves the best results across all scripts. Among open-source models, ViMUL performs favorably against all existing methods across all language scripts, except Chinese, with accuracy ranging from 32% to 55%. Additionally, all video LMMs struggle with *Sinhalese* (Sinhala) and *Tamil* scripts, with ViMUL outperforming the second-best open-source model by 8%. A performance gradient is observed, with video LMM performing significantly better on *Latin*, *Chinese*, and *Cyrillic* scripts compared to *Sinhalese* (Sinhala) and *Tamil*.

We perform an error analysis on cultural examples from ViMUL-Bench by selecting one high-resource language (*Bengali*) and one low-resource language (*Sinhala*), representing the Bengali and Sinhalese scripts, respectively. Native speakers review the open-ended subset responses within the cultural category generated by GPT-4o. Errors are categorized into seven types: lack of knowledge, lack of cultural understanding, language errors, reasoning errors, perceptual errors, translation errors (Vayani et al., 2025; Yue et al., 2024), and prior-

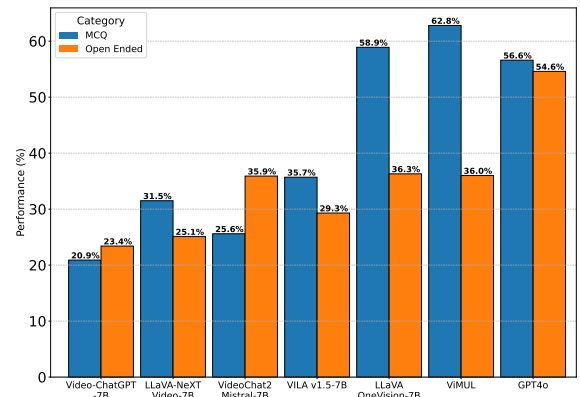


Figure 7: **Performance of different question types in ViMUL-Bench.** Overall, the MCQs attain better performance than the open-ended QAs. ViMUL shows competitive results among open-source models.

knowledge bias. We introduce prior-knowledge bias as a new error type, where the model uses prior knowledge to answer a question with information not present in the video.

Fig. 21 (suppl. material) summarizes the distribution of these error types across *Bengali* and *Sinhalese* scripts, showing that the main errors are knowledge gaps, cultural understanding gaps, prior-knowledge bias, and lack of reasoning capabilities. Fig. 6 presents error examples from language scripts. For instance, in *Bengali*, a question about goal type is asked where GPT-4o incorrectly predicts it as a corner kick, despite the video showing a short throw-in and a header.

**Performance Comparison across Categories.** We study the generic and cultural understanding on 15 categories (Fig. 2b). Overall, GPT-4o achieves best results of 55.8%, but its accuracy varies significantly across different domains. For instance, it scores 75.14% on Festivals but drops to 49.76% on Artistic Performance. In contrast, categories such

as *Digital Content*, *Knowledge*, *Sports Competitions*, and *Surveillance* require a deeper understanding of visual content, likely leading to lower overall performance. Notably, ViMUL surpasses GPT-4o in Media & Entertainment, achieving 57.68%. Similarly, ViMUL outperforms GPT-4o in Notable Public Figures (56.98% vs. 55.87% for GPT-4o) and Sports (70.23% vs. 68.01% for GPT-4o). ViMUL achieves favorable performance on different categories: It achieves 43.29% in *Knowledge* and 45.78% in *Digital Content*, thereby being competitive with GPT-4o. These results show that ViMUL serves as a strong baseline for multilingual, culturally-diverse video understanding.

**Assessing the Need for a Multilingual Video Benchmark.** To motivate the design of ViMUL-Bench, we conduct three baseline ablations demonstrating the necessity of multilingual video input for fair model evaluation. (1) *Blind Baseline*: We show that removing the visual input significantly degrades performance, highlighting the importance of input videos for a fair assessment. (2) *Image-Only Baseline*: LMMs, when evaluated using only single frames (*first/middle/last/random*), exhibit substantial performance drops compared to evaluations using the full video (32 frames), indicating that image-based LMMs are insufficient for capturing the spatio-temporal dynamics required in video benchmarks. (3) *Performance on a Controlled Benchmark*: We assess the impact of our multilingual ViMUL-Instruct fine-tuning on the controlled CVRR-ES (Khattak et al., 2024) benchmark. Models show notable gains in spatio-temporal understanding, particularly in categories requiring cultural or social reasoning. Full experimental results and visualizations are provided in the suppl. material (Sec. L, Fig. 18, Fig. 19, Tab. 3).

**Impact of Location-aware Information in Prompts:** Tab. 2 (suppl. material) presents the comparison when provided with additional country-specific information. Results are consistent with observations in Vayani et al. (2025), with improvement due to better utilization of geographic context.

**Impact of Video Duration.** We further group the videos in ViMUL-Bench into three broad categories based on their duration: short, medium, and long, and present our results in Fig. 20 (suppl. material). Overall, GPT-4o outperforms other models on short and medium videos. However, ViMUL surpasses GPT-4o and other methods on long videos in the multilingual setting. Further details are provided in Sec. K (suppl. material).

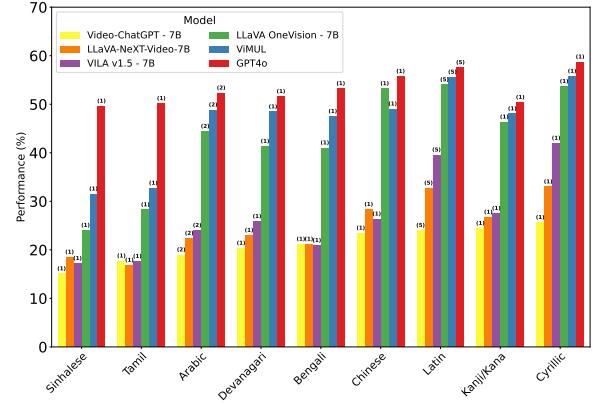


Figure 8: **Performance on different language scripts in ViMUL-Bench.** Models fare higher on high-resource language scripts, such as Latin and Cyrillic, and Chinese, but struggle with under-represented language scripts, such as Tamil and Sinhalese.

## 6 Conclusion

We introduce ViMUL-Bench, the first multilingual benchmark explicitly designed to evaluate video LMMs across diverse linguistic and cultural scenarios. It comprises over 8k humanly verified QA pairs across 14 languages, including both high-resource and several low-resource languages and spanning 15 diverse categories. Further, we present a large-scale multilingual video instruction tuning dataset comprising 1.2 million samples, which we use to develop a simple multilingual video LMM demonstrating competitive cross-linguistic and cultural comprehension.

## Acknowledgement

The computations were enabled through resources provided by NAISS at Alvis, partially funded by the Swedish Research Council under grant agreement no. 2022-06725, LUMI, hosted by CSC (Finland) and the LUMI consortium, and the Berzelius resource provided by the Knut and Alice Wallenberg Foundation at the NSC. This work was partially supported by the Swedish Research Council (2022-04266), from KAW (DarkTree project; 2024.0076) and VR starting grant (2016-05543).

## 7 Limitations

ViMUL, to the best of our knowledge, is the first multilingual Video LMM benchmark that exhibits cultural and linguistic inclusivity across 14 languages and 15 diverse domains. However, despite being the first benchmark for evaluating Video LMMs in a multilingual setting, it has certain limitations. ViMUL-Bench contains more general VQA samples in its sub-categories than in the cultural section, largely due to the high cost of curating culturally grounded data, which demands extensive human effort. It also has more short and medium-length videos than long ones, as high-quality long videos are harder to source for low-resource languages such as Sinhala, Tamil, and Urdu, and require extra verification for reliable QA pairs. While we applied cycle consistency checks during multilingual instruction fine-tuning, further human validation could improve quality. Finally, the ViMUL model, though a simple baseline, highlights the trade-off between high- and low-resource languages in video understanding. As future research directions, our work can be extended to include additional languages, particularly those with limited digital representation. Furthermore, developing culturally-specific, large-scale training datasets tailored explicitly to underrepresented communities would further enhance the inclusivity and effectiveness of multilingual video LMMs. We also expect culture-specific video-language preference data collection can help improve LMMs' performance further via RL.

## 8 Ethical Consideration

Our work reports a standardized multilingual video LMM evaluation benchmark. We hope both ViMUL and ViMUL-Bench will contribute to more consistent evaluation across diverse domains, particularly for underrepresented languages in VidLMM research. Since the cultural videos in ViMUL-Bench are sourced from Internet, some domains may be under-represented, leading to potential biases. To ensure high translation quality, fluent native speakers thoroughly reviewed and verified our initial GPT-4o-generated translations for consistency and accuracy. The verification involves 16 volunteers from diverse linguistic backgrounds, requiring having familiarity with the cultural context of the specific country-language pair they worked on. Additional annotator demographic details are presented in Sec. D suppl. material.

## References

- Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J Hewett, Mojan Javaheripi, Piero Kauffmann, and 1 others. 2024. Phi-4 technical report. *arXiv preprint arXiv:2412.08905*.
- Meta AI. 2024. Llama 3. <https://llama.meta.com/llama3>.
- Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, and 1 others. 2023. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *arXiv preprint arXiv:2302.04023*.
- Lyle Campbell and Verónica Grondona. 2008. Ethnologue: Languages of the world. *Language*, 84(3):636–641.
- Ron Campos, Ashmal Vayani, Parth Parag Kulkarni, Rohit Gupta, Aritra Dutta, and Mubarak Shah. 2025. Gaea: A geolocation aware conversational model. *arXiv preprint arXiv:2503.16423*.
- Marta R Costa-Jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, and 1 others. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.
- Rocktim Jyoti Das, Simeon Emilov Hristov, Haonan Li, Dimitar Iliyanov Dimitrov, Ivan Koychev, and Preslav Nakov. 2024. Exams-v: A multi-discipline multilingual multimodal exam benchmark for evaluating vision language models. *arXiv preprint arXiv:2403.10378*.
- Bernard Ghanem Fabian Caba Heilbron, Victor Escorcia and Juan Carlos Niebles. 2015. Activitynet: A large-scale video benchmark for human activity understanding.
- Chaoyou Fu, Yuhang Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, and 1 others. 2024. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. *arXiv preprint arXiv:2405.21075*.
- C Haerpfer, R Inglehart, A Moreno, C Welzel, K Kizilova, J Diez-Medrano, M Lagos, P Norris, E Ponarin, and B Puranen. 2022. World values survey: Round seven-country-pooled datafile version 3.0. jd systems institute & wvsa secretariat.
- Harald Hammarström, Robert Forkel, Martin Haspelmath, and Sebastian Bank. 2022. *Glottolog database* 4.6.
- Weiwei Hu, Yunji Chen, Tianshi Chen, Cheng Qian, and Lei Li. 2011. Linear time memory consistency verification. *IEEE Transactions on Computers*, 61(4):502–516.

- Ilker Kesen, Andrea Pedrotti, Mustafa Dogan, Michele Cafagna, Emre Can Acikgoz, Letitia Parcalabescu, Iacer Calixto, Anette Frank, Albert Gatt, Aykut Erdem, and Erkut Erdem. 2024. Vilma: A zero-shot benchmark for linguistic and temporal grounding in video-language models. In *International Conference on Learning Representations (ICLR)*.
- Muhammad Uzair Khattak, Muhammad Ferjad Naeem, Jameel Hassan, Muzammal Naseer, Federico Tomba, Fahad Shahbaz Khan, and Salman Khan. 2024. How good is my video lmm? complex video reasoning and robustness evaluation suite for video-lmms. *arXiv preprint arXiv:2405.03690*.
- Jie Lei, Licheng Yu, Mohit Bansal, and Tamara L Berg. 2018. Tvqa: Localized, compositional video question answering. *arXiv preprint arXiv:1809.01696*.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and 1 others. 2024a. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*.
- Bohao Li, Yuying Ge, Yixiao Ge, Guangzhi Wang, Rui Wang, Ruimao Zhang, and Ying Shan. 2023a. Seed-bench-2: Benchmarking multimodal large language models. *arXiv preprint arXiv:2311.17092*.
- Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. 2023b. Seed-bench: Benchmarking multimodal llms with generative comprehension. *arXiv preprint arXiv:2307.16125*.
- KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. 2023c. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*.
- Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, and 1 others. 2024b. Mvbench: A comprehensive multi-modal video understanding benchmark. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22195–22206.
- Fangyu Liu, Emanuele Bugliarello, Edoardo Ponti, Siva Reddy, Nigel Collier, and Desmond Elliott. Visually grounded reasoning across languages and cultures. In *NeurIPS 2021 Workshop on ImageNet: Past, Present, and Future*.
- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, and 1 others. 2024. Mmbench: Is your multi-modal model an all-around player? In *European conference on computer vision*, pages 216–233. Springer.
- Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. 2024a. Video-chatgpt: Towards detailed video understanding via large vision and language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL 2024)*.
- Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. 2024b. Videogpt+: Integrating image and video encoders for enhanced video understanding. *arxiv*.
- Muhammad Maaz, Hanoona Rasheed, Abdelrahman Shaker, Salman Khan, Hisham Cholakkal, Rao M Anwer, Tim Baldwin, Michael Felsberg, and Fahad S Khan. 2025. Palo: A polyglot large multimodal model for 5b people. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*.
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*, pages 3195–3204.
- Vishal Narnaware, Ashmal Vayani, Rohit Gupta, Swetha Sirnam, and Mubarak Shah. 2025. Sb-bench: Stereotype bias benchmark for large multimodal models. *arXiv preprint arXiv:2502.08779*.
- OpenAI. 2024. Gpt-4v. <https://api.semanticscholar.org/CorpusID:263218031>.
- Viorica Patraucean, Lucas Smaira, Ankush Gupta, Adria Recasens, Larisa Markeeva, Dylan Banarse, Skanda Koppula, Mateusz Malinowski, Yi Yang, Carl Doersch, and 1 others. 2023. Perception test: A diagnostic benchmark for multimodal video models.
- Jonas Pfeiffer, Gregor Geigle, Aishwarya Kamath, Jan-Martin O Steitz, Stefan Roth, Ivan Vulić, and Iryna Gurevych. 2021. xgqa: Cross-lingual visual question answering. *arXiv preprint arXiv:2109.06082*.
- Rizwan Qureshi, Ranjan Sapkota, Abbas Shah, Amgad Muneer, Anas Zafar, Ashmal Vayani, Maged Shoman, Abdelrahman Eldaly, Kai Zhang, Ferhat Sadak, and 1 others. 2025. Thinking beyond tokens: From brain-inspired intelligence to cognitive foundations for artificial general intelligence and its societal impact. *arXiv preprint arXiv:2507.00951*.
- Shaina Raza, Rizwan Qureshi, Anam Zahid, Joseph Fiorelli, Ferhat Sadak, Muhammaed Saeed, Ranjan Sapkota, Aditya Jain, Anas Zafar, Muneeb Ul Hassan, and 1 others. 2025a. Who is responsible? the data, models, users or regulations? responsible generative ai for a sustainable future. *Authorea Preprints*.
- Shaina Raza, Ashmal Vayani, Aditya Jain, Aravind Narayanan, Vahid Reza Khazaie, Syed Raza Bashir, Elham Dolatabadi, Gias Uddin, Christos Emmanouilidis, Rizwan Qureshi, and 1 others. 2025b. Vld-bench: Vision language models disinformation detection benchmark. *arXiv e-prints*, pages arXiv–2502.
- David Romero, Chenyang Lyu, Haryo Akbarianto Wibowo, Teresa Lynn, Injy Hamed, Aditya Nanda Kishore, Aishik Mandal, Alina Dragonetti, Artem Abzaliev, Atnafu Lambebo Tonja, and 1 others.

2024. Cvqa: Culturally-diverse multilingual visual question answering benchmark. *arXiv preprint arXiv:2406.05967*.
- Chenhui Shen, Liying Cheng, Xuan-Phi Nguyen, Yang You, and Lidong Bing. 2023. Large language models are not yet human-level evaluators for abstractive summarization. *arXiv preprint arXiv:2305.13091*.
- Rickard Stureborg, Dimitris Alikaniotis, and Yoshi Suhara. 2024. Large language models are inconsistent and biased evaluators. *arXiv preprint arXiv:2405.01724*.
- Hai-Long Sun, Da-Wei Zhou, Yang Li, Shiyin Lu, Chao Yi, Qing-Guo Chen, Zhao Xu, Weihua Luo, Kaifu Zhang, De-Chuan Zhan, and 1 others. 2024. Parrot: Multilingual visual instruction tuning. *arXiv preprint arXiv:2406.02539*.
- Makarand Tapaswi, Yukun Zhu, Rainer Stiefelhausen, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. 2016. Movieqa: Understanding stories in movies through question-answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4631–4640.
- Qwen Team. 2024. [Hello qwen2](#).
- Ashmal Vayani, Dinura Dissanayake, Hasindri Watawana, Noor Ahsan, Nevasini Sasikumar, Omkar Thawakar, Henok Biadglign Ademtew, Yahya Hmaiti, Amandeep Kumar, Kartik Kuckreja, and 1 others. 2025. All languages matter: Evaluating llms on culturally diverse 100 languages. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Hongyu Wang, Jiayu Xu, Senwei Xie, Ruiping Wang, Jialin Li, Zhaojie Xie, Bin Zhang, Chuyan Xiong, and Xilin Chen. 2024. M4u: Evaluating multilingual understanding and reasoning for large multimodal models. *arXiv preprint arXiv:2405.15638*.
- Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. 2021. Next-qa: Next phase of question-answering to explaining temporal actions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9777–9786.
- Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. 2017. Video question answering via gradually refined attention over appearance and motion. In *Proceedings of the 25th ACM international conference on Multimedia*, pages 1645–1653.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.
- Kexin Yi, Chuang Gan, Yunzhu Li, Pushmeet Kohli, Jiajun Wu, Antonio Torralba, and Joshua B Tenenbaum. 2019. Clevrer: Collision events for video representation and reasoning. *arXiv preprint arXiv:1910.01442*.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, and 1 others. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567.
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. Sigmoid loss for language image pre-training.
- Hongjie Zhang, Yi Liu, Lu Dong, Yifei Huang, Zhen-Hua Ling, Yali Wang, Limin Wang, and Yu Qiao. 2023a. Movqa: A benchmark of versatile question-answering for long-form movie understanding. *arXiv preprint arXiv:2312.04817*.
- Wenxuan Zhang, Mahani Aljunied, Chang Gao, Yew Ken Chia, and Lidong Bing. 2023b. M3exam: A multilingual, multimodal, multilevel benchmark for examining large language models. *Advances in Neural Information Processing Systems*, 36:5484–5505.
- Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. 2024. Video instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*.
- Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. 2016. Visual7w: Grounded question answering in images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4995–5004.

## A Prompts used in Evaluation

**Multiple-Choice Evaluation Prompt.** All LMMs are prompted on the below multiple-choice evaluation prompt to select the most accurate answer from the available options based on both the video sample and the subtitles. The model is strictly instructed to select only one correct option (A, B, C, or D), ensuring a clear and direct response format. Then, the model's generated output is compared directly with the ground-truth to evaluate accuracy. The prompt for MCQ evaluation is shown below:

### LMM Inference Prompt.

**Prompt:**

Select the best answer to the following multiple-choice question based on the video and the subtitles. Respond with only the letter (A, B, C, or D) of the correct option.

**Open-Ended Evaluation Prompt.** Below, we list the evaluation guidelines for the open-ended evaluation.

### Evaluation Guidelines

**System Prompt:**

You are an intelligent chatbot designed for evaluating the correctness of AI assistant predictions for question-answer pairs. Your task is to compare the predicted answer with the ground-truth answer and determine if the predicted answer is correct or not. Here's how you can accomplish the task:

**Instructions:**

- **Correctness Focus:** Compare the predicted answer with the ground-truth answer to determine accuracy.
- **Detail Consideration:** Predictions with fewer details are still correct unless such details are explicitly required in the question.
- **Scoring:** Assign an integer score between 0 (fully incorrect) and 5 (fully correct), with intermediate values reflecting partial correctness.

This GPT prompt is designed for evaluating AI assistant predictions on video-based question-

answer pairs. It instructs the model to compare predicted responses with ground-truth answers, assessing correctness while allowing for minor variations unless explicitly required. The evaluation follows a structured scoring system from 0 (fully incorrect) to 5 (fully correct), with intermediate scores reflecting partial accuracy. The output is formatted as a Python dictionary containing the prediction status ("correct" or "incorrect"), a numerical score, and a justification, ensuring consistency in automated assessment for multilingual and multimodal AI systems.

### Evaluation Request.

Please evaluate the following video-based question-answer pair:

**Question:** {question}

**Ground Truth Correct Answer:** {answer}

**Predicted Answer:** {pred}

Provide your evaluation as a correct/incorrect prediction along with the score, which is an integer value between 0 (fully wrong) and 5 (fully correct). The middle score represents the percentage of correctness.

**Response Format:**

Your response should be generated as a Python dictionary string with the following keys:

- **'pred':** A string, either "correct" or "incorrect".
- **'score':** An integer between 0 and 5.
- **'reason':** A justification for the decision.

Only provide the Python dictionary string.

Example format:

```
{"pred": "correct", "score": 4,
"reason": "The predicted answer captures most of the ground-truth meaning but lacks minor details."}
```

## B Dataset Statistics

**Cultural vs Generic Distribution.** ViMUL-Bench includes a total of 15 diverse categories, comprising both generic and cultural categories. Figure 10 illustrates the distribution of these cate-




French: Latin Script	Japanese: Kanji/ Kana Script	Spanish: Latin Script
 <p><b>Question:</b> Quels drapeaux apparaissent dans la vidéo ?  <b>Ground Truth Answer:</b> France et Union européenne  <b>Predicted Answer:</b> Dans la vidéo, les drapeaux visibles sont le drapeau français (bleu, blanc, rouge) et le drapeau de l'Union européenne (bleu avec des étoiles jaunes).  <b>Score:</b> 5/5  <b>Video Duration:</b> Medium (6+ Mins)</p>	 <p><b>Question:</b> 公開された屏風に描かれているのは、何の物語ですか？(A. 源氏物語, B. 日本書紀, C. 竹取物語, D. 伊勢物語)  <b>Ground Truth Answer:</b> 源氏物語  <b>Predicted Answer:</b> A  <b>Score:</b> Accurate  <b>Video Duration:</b> Medium (9+ Mins)</p>	 <p><b>Question:</b> ¿Qué artista aparece en el video? (A. Vincent van Gogh, B. Salvador Dali, C. Frida Kahlo, D. Pablo Picasso)  <b>Ground Truth Answer:</b> Salvador Dali  <b>Predicted Answer:</b> B  <b>Score:</b> Accurate  <b>Video Duration:</b> Short (3+ Mins)</p>
<b>Category:</b> Notable Public Figures	<b>Category:</b> Architecture & Landmarks	<b>Category:</b> Art & Literature

Figure 9: We present qualitative examples of success cases of GPT-4o’s across different language scripts and categories.

gories across different sources. The generic samples are carefully selected from VCG-Diverse, CVRR-ES, MVBench, and VideoMME datasets. For the cultural category, we curate the content from scratch with input from native speakers for all languages except English. Additionally, we use GPT-4o to translate video question-answer pairs from all 13 languages into English.

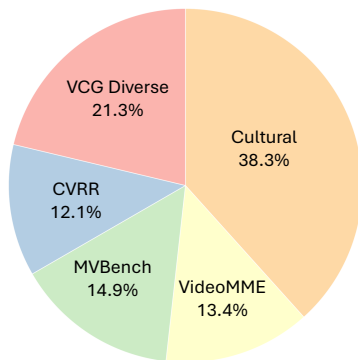


Figure 10: We present the cultural versus generic category distribution. We source generic categories from existing video benchmarks. However, the cultural part is carefully curated from scratch by native speakers.

**Video Duration Distribution.** ViMUL-Bench categorizes videos into three durations: short (0-4 minutes), medium (4-15 minutes), and long (15+ minutes). Each video is manually assigned a duration label. Figure 11 illustrates the distribution of video durations across the dataset. The plot shows that over 73.7% of the videos are short, making it the largest category, followed by 20.5% medium-duration videos and 5.8% long-duration videos. Although we aimed to achieve a balanced distribution of video durations, curating long-duration videos proved to be more resource-intensive for verification, and they are relatively scarce.

**Language Distribution.** Figure 12 illustrates the language distribution across ViMUL-Bench. The

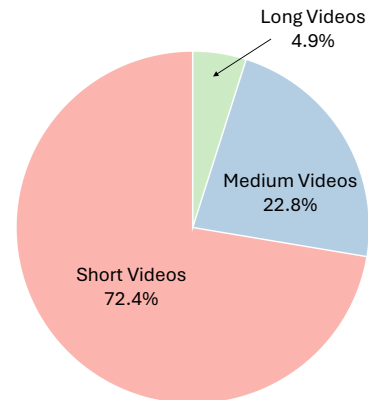


Figure 11: The figure illustrates the video duration distribution in ViMUL-Bench. Overall, we have over 73% short-duration videos in the dataset.

dataset contains a nearly equal proportion of videos for each of the 14 languages, with the exception of English, which is more heavily represented. For the cultural data, we translate video question-answer pairs from all 13 non-English languages into English to ensure consistency and facilitate cross-lingual comparisons. This approach enables us to maintain linguistic diversity while ensuring accessibility and usability of the dataset in English.

**ViMUL Category Distribution.** Our benchmark consists of 15 diverse categories, including seven generic categories and eight cultural categories. Figure 14 presents the distribution of the seven generic categories. Among these, the Life Record category accounts for approximately 23% of the samples, followed by Digital Content (17.5%) and Knowledge (11.1%). Surveillance and Artistic Performance categories represent smaller proportions of the dataset. In total, over 33% of the dataset is composed of cultural categories, which are further divided into eight subcategories. Fig. 13 shows the cultural category distribution, where all the categories are almost balanced, ranging from 15.1%

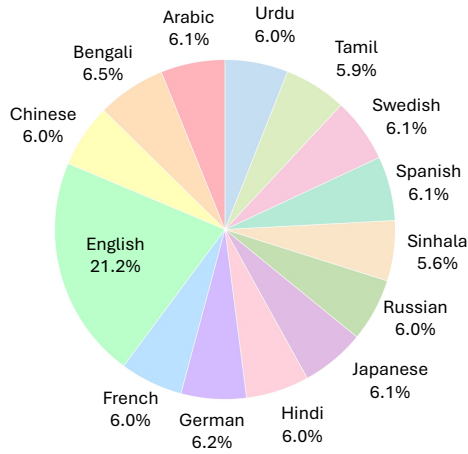


Figure 12: The figure illustrates the per-language distribution in ViMUL-Bench. The dataset contains nearly equal proportions of both low-resource and high-resource languages, except English. English comprises translations of video question-answer pairs from all 13 other languages in the dataset.

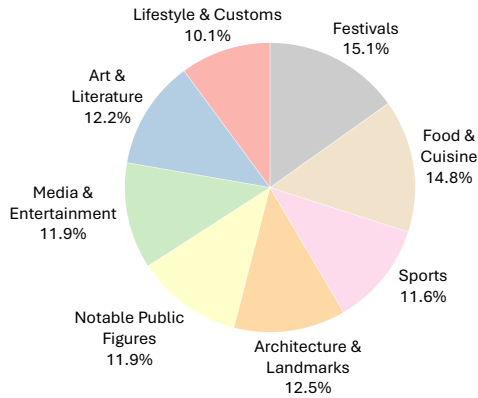


Figure 13: The figure illustrates the Distribution of the eight cultural categories in ViMUL-Bench, where we ensure consistent samples across all the cultural categories.

of QA pairs in Festival to 10.1% in Lifestyle and Customs.

### C Impact of Location-aware Information in Prompts

Table 2 presents the performance improvements when geographical information, specifically country details, is included in the prompts. The closed-source model demonstrates a greater ability to leverage this additional information compared to the open-source model.

### D Volunteer Demographics

We have a total of 16 volunteers from various backgrounds who assisted us in curating and verify-

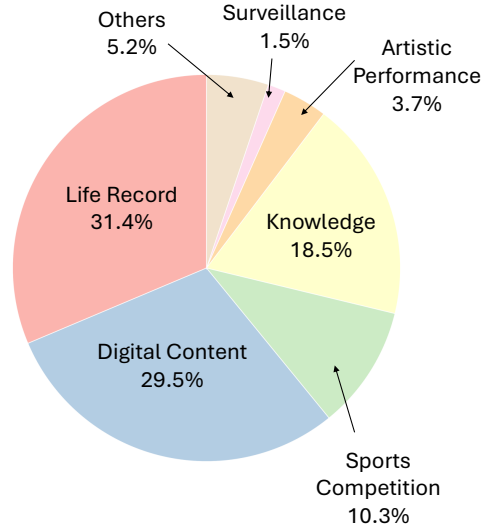


Figure 14: The figure illustrates the Distribution of the seven generic categories in ViMUL-Bench, with the largest proportions in Life Record, Digital Content, and Knowledge from generic categories.

Models	With Country Info.	Without Country Info.
GPT-4o	63.3%	60.8%
ViMUL	53.1%	52.61%

Table 2: **Performance with and without additional country location information.** Results improve when integrating additional geographic information as input to VidLMs.

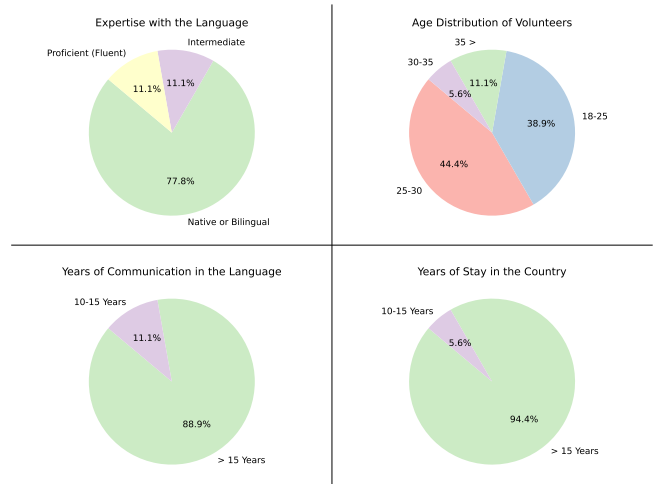


Figure 15: The top left figure shows the percentage of our volunteers with respect to linguistic skill level. The top right shows their age distribution. Then in bottom left we have their active years of communication, followed by the duration of their stay in their respective countries.

ing our ViMUL Bench. Among these, two volunteers are proficient in two different languages each,

and they contributed separately for both languages. To accurately reflect their contributions and language proficiency levels, we consider a total of 18 language instances rather than just 16 individuals when calculating statistics.

Around 77.8% of our verifiers are native or bilingual, followed by 11.1% of them proficient and the rest are intermediate in their respective languages. Around 94.4% stayed for more than 15 years in their country, where they learned their first language. We can observe in Fig. 15 our contributors range from the age bracket of 18 - 25, 25 - 30 and onwards, which makes the age distribution more diverse, and includes prior experience. In terms of geographical distribution, they come from the following countries as follows Bangladesh, China, Germany, India, Japan, Lebanon, Morocco, Pakistan, Sri Lanka, Sweden, Ukraine and USA. This diversity in geo-location helps us to cater to cultural nuances of that corresponding region, hence allowing us to get translations and verifications which are authentic to that place.

## E Cycle Consistency for Train Set

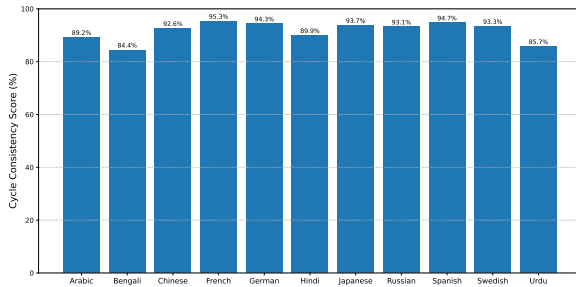


Figure 16: **Performance of Cycle consistency upon 11 languages.** All the 11 languages display an accuracy well above 80%, with the average being 91.47% across the languages.

We perform a cycle consistency check to evaluate the quality of machine translations across multiple languages. Starting with a pre-existing English-only training dataset, we translate the samples into 13 additional languages using GPT-4o-mini. To assess the accuracy and fidelity of these translations, we then back-translate the non-English samples into English using Qwen 3. The back-translations are evaluated based on the following criteria: consistency with the original text, grammatical correctness, and meaningfulness. These evaluations are conducted using GPT-4o, allowing us to systematically score and compare machine translation in terms of both linguistic and semantic properties

between languages. We can observe in Fig. 16 that the highest performance can be seen in French with 95.3%, followed by Spanish with 94.7%. At the lower end, we can find Bengali with 84.4%, followed by Urdu with 85.7%.

## F Performance of Phi-4 as a judge

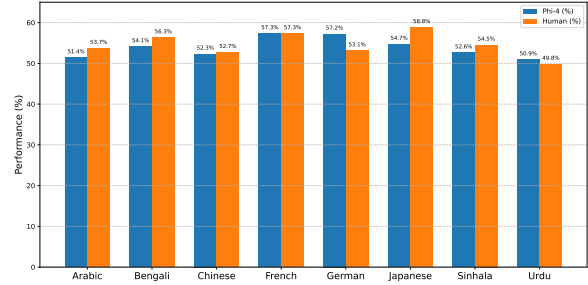


Figure 17: **Performance of Phi-4 scores compared to human score.** The figure illustrates that Phi-4 scores are similar to human scores where they differ most in Japanese by 4.1% and the least in Chinese by 0.4%.

To validate the reliability of Phi-4 as a judge, we conducted human verification on 100 randomly sampled open-ended VQAs scores (both cultural and generic categories) from our ViMUL-Bench. Fig. 17 (suppl.) shows the performance consistency on GPT-4o responses, where all languages, both low resource and high resource languages, exhibit strong agreement with Phi-4’s scores when compared to manual native speaker scores. For instance, for Urdu language, the average accuracy of Phi-4 score is 50.9%, whereas for human score the average is 49.8%. Hence we can deduce that since the Phi-4 scores are similar to human scores, thus Phi-4 scores are reliable.

## G Performance of model based upon question only

Fig. 18 compares the performance of ViMUL, LLaVA OneVision, and LLaVA-NeXT across 14 languages when only the question is provided as input (no frames). The results highlight that ViMUL achieves a slightly higher average performance of 28.5%, outperforming LLaVA OneVision, which has an average of 28.3%, and LLaVA-NeXT, which averages 22.7%. ViMUL demonstrates a better overall tradeoff between high-resource and low-resource languages. Specifically, in languages such as Sinhala and Urdu, which are considered low-resource, ViMUL maintains relatively stronger performance compared to LLaVA OneVision and

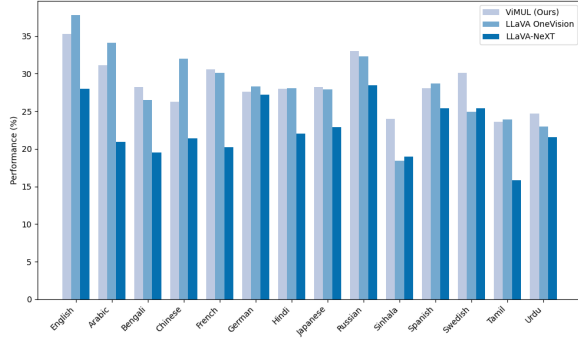


Figure 18: **Performance of models for question only as input.** This Figure shows the performance of models for all 14 languages when only question is given as input [No frames].

LLaVA-NeXT. This shows that ViMUL is more robust across both high-resource languages (like Russian, French and Swedish) and low-resource languages, making it a more versatile model for multilingual tasks.

## H Performance of ViMUL-LLM before and after finetuning

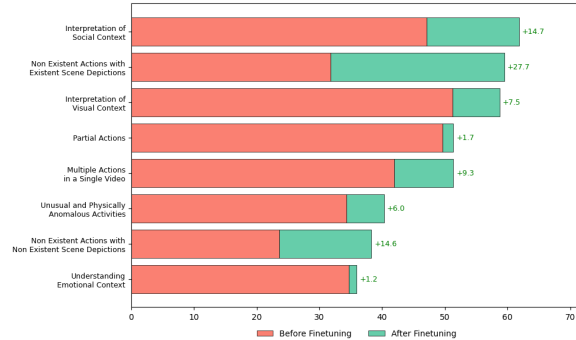


Figure 19: **Performance before Vs after finetuning.** This Figure portrays the performance enhancement of ViMUL model after finetuning for various tasks, such as spatial and temporal.

Fig. 19 demonstrates the performance of the ViMUL model before and after fine-tuning across eight distinct tasks. The tasks include a variety of challenges such as interpreting social and visual contexts, handling multiple actions in a single video, and recognizing emotionally charged or anomalous activities. The figure shows the improvement in performance after fine-tuning, with the score for each task before and after fine-tuning represented by horizontal bars. The difference in scores, indicated by the green bars, highlights the model’s enhancement across different tasks, particularly in areas like *Interpretation of Social Context*

and recognizing *Non Existent Actions with Existent Scene Depictions*, which increases by 14.7% and 27.7%, respectively. This visual comparison underscores the effectiveness of the fine-tuning process in boosting the overall accuracy of the model and the specific performance of the task.

## I Performance based upon frames

Setting	LLaVA-OneVision	ViMUL (Ours)
First Frame	34.1	35.3
Mid Frame	36.7	39.0
Last Frame	34.6	35.8
Random Frame	37.7	38.4
ViMUL-Bench	46.2	49.2

Table 3: Comparison of image and video backbones. Accuracy of LLaVA-NeXT, LLaVA OneVision, and ViMUL across different frame inputs.

Tab. 3 compares the performance of LLaVA-OneVision and ViMUL model across different input types, including first, mid, last, and random frames, as well as when 32 frames are taken as input, referred to as "ViMUL-Bench." The results demonstrate that ViMUL consistently outperforms LLaVA-OneVision across all frame-based settings. Notably, ViMUL achieves the highest accuracy of 49.2% in the ViMUL-Bench setting, where 32 frames are used as input. This significant performance boost underscores the importance of videos over individual images. Videos provide crucial temporal context that images alone cannot capture, allowing models to leverage the dynamics and sequences of actions, resulting in better overall accuracy. The improvement observed with multiple frames emphasizes that incorporating temporal information in video inputs is key to enhancing model performance.

## J ViMUL: Multilingual Video LMM

**Video Sampling:** Given an input video  $\mathbf{V} \in \mathbb{R}^{T \times H \times W \times C}$ , where  $T$  is the total number of frames, and  $H$  and  $W$  denote the height and width of each frame respectively, we define  $N$  as the maximum number of frames that can be processed. We first sample the video at 1 FPS, resulting in  $t$  frames. If  $t \leq N$ , we keep all sampled frames. However, if  $t > N$ , we further uniformly select  $N$  frames from the  $t$  frames. This ensures that the final number of frames does not exceed  $N$ . Thus, the video representation after sampling is  $\mathbf{V}' \in \mathbb{R}^{n \times H \times W \times C}$ ,

where  $n = \min(N, t)$  and  $\mathbf{V}'$  represents the sampled video with at most  $N$  frames.

**Video Encoding:** We use SigLIP (Zhai et al., 2023) as the vision encoder, followed by a vision-to-language projector to transform vision tokens into the input embedding space of the language model. Given the video  $\mathbf{V}' \in \mathbb{R}^{n \times H \times W \times C}$ , we first resize the input to the resolution that SigLIP is trained on, followed by encoding each frame of the video using SigLIP. The output from the second-to-last layer of SigLIP is flattened and passed through a two-layer MLP, which projects the vision features into the language model’s embedding space. The projected features are reshaped into a grid format and pooled using a  $2 \times 2$  kernel to reduce the number of features by a factor of four. Empirically, we found that this helps accommodate more video frames while maintaining performance. Finally, the projected features are flattened, resulting in video embeddings  $\mathbf{E}^{vid} \in \mathbb{R}^{n \times L_v \times D_t}$ , where  $L_v$  represents the total visual features per video frame, and  $D_t$  is the embedding size of the language model.

**Language Model:** We obtain the final representation by concatenating the video embeddings  $\mathbf{E}^{vid}$  with the text embeddings  $\mathbf{E}^{text} \in \mathbb{R}^{L \times D_t}$  of the user query,

$$\mathbf{E} = [\mathbf{E}^{vid}, \mathbf{E}^{text}]. \quad (1)$$

This ensures the language model receives spatio-temporal video features followed by the user query to generate an accurate response. We use Qwen-2.0 (Team, 2024) as the language model and fully fine-tune it in an auto-regressive manner with a next-token prediction loss (see Fig. 4).

**Dataset Distribution:** The LLaVA-Video-178K dataset includes 178,510 caption entries, 960,792 open-ended QA pairs, and 196,198 multiple-choice samples. We employ LLaMA-3.1-70B-Instruct (AI, 2024) model to identify the complex QA pairs in both open-ended and multiple-choice categories. Specifically, we prompt the LLM to identify samples that are difficult and require chain-of-thought reasoning, resulting in 39,422 QA pairs. Further, we include the training sets of NeXT-QA (Xiao et al., 2021), PerceptionTest (Patraucean et al., 2023), and Clevrer (Yi et al., 2019), contributing an additional 29,846 samples. This results in a total of 95k samples in our English training dataset.

We translate these English QA pairs into 13 other languages, including Arabic, Bengali, Chinese, French, German, Hindi, Japanese, Russian,

Sinhala, Spanish, Swedish, Tamil, and Urdu, using GPT-4o-mini (OpenAI, 2024). We observe that GPT-4o occasionally makes obvious mistakes in translation, such as failing to respond in the target language. To address these issues, we employ LLaMA-3.1-8B (AI, 2024) model to post-process the translations. We input the translated QA pairs to the LLM and ask it to predict the language in which the text is written. If the LLM predicts that the text is not in the intended language, we simply discard these samples.

Finally, our pipeline results in a total of 1,238,102 samples across all languages. Original English dataset contains 95,071 samples, while translated datasets contain; Arabic: 88,154, Bengali: 88,087, Chinese: 88,095, French: 86,990, German: 88,020, Hindi: 88,004, Japanese: 88,041, Russian: 88,054, Sinhala: 88,023, Spanish: 87,976, Swedish: 87,974, Tamil: 87,946, Urdu: 87,667 samples, respectively.

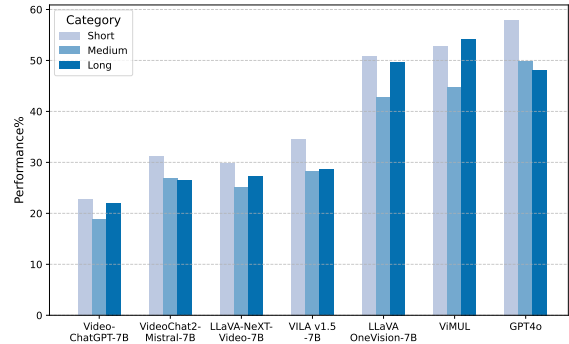


Figure 20: **Performance comparison across three video durations in ViMUL-Bench.** Most methods perform better on short video questions, followed by long and then medium-length videos. GPT-4o achieves the highest accuracy on short VQA, while ViMUL outperforms all LMMs on long VQAs.

## K Impact of Video Duration

As mentioned earlier, video samples in ViMUL-Bench are grouped into three categories based on their duration: short (0-4 mins), medium (4-15 mins), and long (15+ mins). Fig. 20 illustrates how different methods perform across these categories. Overall, GPT-4o achieves the highest accuracy on short and medium videos but struggles with long videos in a multilingual setting. In comparison, our proposed baseline, ViMUL, outperforms all LMMs on long videos and surpasses LLaVA-OneVision on both short and medium videos. Most open-source models perform best on short videos, with their accuracy struggling on medium and long videos.

The only exception is ViMUL, due to its extensive multilingual training corpus. Notably, the difference between short and long video performance is comparable in many open-source methods, possibly due to the fewer long videos as compared to short videos in ViMUL-Bench, as shown in Fig. 13 (suppl. material). The culturally curated long videos are harder to find and expensive to manually annotate.

## L Assessing the need for a Multilingual Video Benchmark

To address this question, we conduct three baseline ablations on ViMUL-Bench. **(1) Blind Baseline:** We evaluate LMMs using only the textual QA pairs, without providing the visual input. As shown in Fig. 18 (suppl. material), the performance of LLM-only variants drops significantly in the absence of video input. This demonstrates the necessity of incorporating visual input to ensure fair and comprehensive model evaluation. **(2) Image LMM Baseline:** We prompt LLaVA-OneVision and ViMUL with single frames instead of the complete video (32 frames). For rigorous evaluation, we test it on the *First*, *Middle*, *Last*, and *Random* frames. As summarized in Tab. 3 (suppl. material), the performance of LLaVA-OneVision drops by 8.9% when prompted with a random frame and by 12% when prompted with the first frame. A similar trend is observed for ViMUL. **(3) Performance on a controlled benchmark:** We conduct an additional experiment using the controlled benchmark CVRR-ES (Khat-tak et al., 2024) to assess model performance before and after fine-tuning on ViMUL-Instruct. We sample eight spatiotemporal dimensions from CVRR-ES and present the results in Fig. 19 (suppl. material). Categories such as *Non-Existent Actions with Existent Scene Depictions* and *Interpretation of Social Context* show improvements of 27.68% and 14.68%, respectively. Other categories, including *Interpretation of Visual Context* and *Multiple Actions in a Single Video*, also demonstrate consistent performance gains. However, categories like *Understanding Emotional Context* and *Partial Actions* exhibit comparatively lower improvement.

## M Error analysis of GPT-4o by native speakers

Fig. 21 displays an error analysis of GPT-4o’s performance by native speakers on ViMUL-Bench across Bengali and Sinhalese scripts, revealing dis-

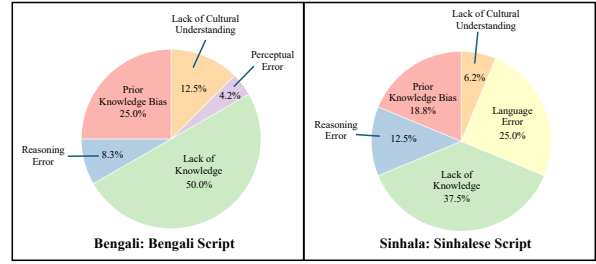


Figure 21: **Performance of different question types in ViMUL-Bench.** Overall, the MCQs attain better performance than the open-ended QAs. ViMUL achieves competitive performance compared to existing open-source models.

tinct patterns of failure. In the Bengali script, the most frequent issues stem from Prior Bias Knowledge and Lack of Knowledge, which are 25% and 50%, respectively. While in the Sinhalese script, errors are more often attributed to Lack of Knowledge and Language Error, which are 37.5% and 25%, respectively. These disparities highlight the complex interplay between language structure, cultural grounding, and model capabilities. ViMUL-Bench effectively surfaces such weaknesses, emphasizing the need for culturally and linguistically grounded evaluation in multilingual VQA systems.

## N License for Artifacts

We use the following datasets and models in our paper. We list all the licenses for each of them below

- Base LLMs and the multimodal finetuned models accessed via Hugging Face (e.g., Video-ChatGPT-7B, LLaVA-NeXT-Video-7B,, VideoChat2-Mistral-7B, VILA v1.5 - 7B, LLaVA OneVision - 7B, GPT-4o, GPT-4o-mini, Qwen-3, LLaMA-3.1, Phi-4) are cited with their respective papers, model cards, and URLs.

Licenses for multimodal LMMs:

- Video-ChatGPT-7B: CC BY 4.0
- LLaVA-NeXT-Video-7B: LLAMA 2 Community License
- VideoChat2-Mistral-7B: Apache License 2.0
- VILA-v1.5-7B: Apache License 2.0 (code), CC BY-NC-SA 4.0 (weights)
- LLaVA-OneVision-7B: LLAMA 2 Community License
- Qwen-3: Apache License 2.0

- Phi-4: MIT License
- GPT-4o: Closed-source (OpenAI Proprietary)
- GPT-4o-mini: Closed-source (OpenAI Proprietary)
- LLaMA-3.1: LLaMA 3 Community License (restricted commercial use)

Licenses for the datasets used in our work for the construction of fine-tuning and evaluation benchmark.

- **CVRR**: CC BY-NC-SA 4.0
- **VideoMME**: Academic research only; commercial use prohibited
- **VCG Diverse**: MIT License
- **MVBench**: MIT License
- **Video-Instruct100K**: CC BY-SA 4.0
- **LLaVA-Video-178K**: Apache License 2.0
- **NeXT-QA**: MIT License
- **PerceptionTest**: CC BY 4.0
- **CLEVRER**: CC0 License
- **ActivityNetQA**: MIT License (non-commercial use only)
- The video samples for cultural category was scraped manually from the publicly accessible websites, such as YouTube. The usage of this content is compliant with fair-dealing law for non-commercial academic research. We do not redistribute the original video and text data under commercial licensing.
- Our codebase builds on open-source frameworks such as PyTorch (BSD-style license) and Transformers (Apache 2.0). Our evaluation framework is based on the Imms-eval framework, based on Apache License 2.0.

The terms of use for each resource are respected, and we provide links and citations in the supplementary material and code repository. We release our work under the license [CC BY-NC 4.0](#).

## O Documentation of Artifacts

We list the coverage of domains of the datasets used in our paper:

### O.1 Evaluation Benchmark Domains

#### • CVRR:

- Non-existent actions with non-existent scene depictions
- Non-existent actions with existent scene depictions
- Time order understanding
- Understanding of emotional context
- Interpretation of social context
- Unusual and Physically Anonymous activities
- Continuity and Object Instance Count
- Interpretation of visual context
- Partial actions
- Fine-grained action understanding
- Multiple actions in a single video

#### • VideoMME:

- Domains in dataset:
  - \* Life Record: Travel, Daily Life, Fashion, Food, Handicraft, Pet & Animal, Exercise, Multilingual
  - \* Artistic Performance: Acrobatics, Variety Show, Magic Show, Stage Play
  - \* Sports Competition: Other Sports, Athletics, Football, Basketball, Esports
  - \* Film & Television: News Report, Documentary, Movie & TV Show, Animation
  - \* Knowledge: Technology, Life Tip, Law, Geography, Astronomy, Finance & Commerce, Biology & Medicine, Literature & Art, Humanity & History
- Domains in various questions:
  - \* Temporal Perception
  - \* Spatial Perception
  - \* Attribute Perception
  - \* Action Recognition
  - \* Object Recognition
  - \* OCR Problems
  - \* Counting Problems
  - \* Temporal Reasoning
  - \* Spatial Reasoning
  - \* Action Reasoning
  - \* Object Reasoning
  - \* Information Synopsis

#### • MVBench:

- Spatial Understanding
  - \* Action: What’s the man doing?
  - \* Object: What’s on the table?
  - \* Position: Is the man on the stage?
  - \* Count: How many chairs?
  - \* Scene: Where’s the man?
  - \* Pose: What’s the man’s pose?
  - \* Attribute: What color is the desk?
  - \* Character: What are the subtitles?
  - \* Cognition: Why is the man singing in the canteen?
- Temporal Understanding
  - \* Action: Action Sequence, Action Antonym, Action Prediction, Unexpected Action, Fine-grained Action
  - \* Object: Object Shuffle, Object Existence, Object Interaction
  - \* Position: Moving Direction, Action Localization
  - \* Count: Action Count, Moving Count
  - \* Scene: Scene Transition
  - \* Pose: Fine-grained Pose
  - \* Attribute: State Change, Moving Attribute
  - \* Character: Character Order
  - \* Cognition: Episodic Reasoning, Ego-centric Navigation, Counterfactual Inference
- **VCG-Diverse**
  - Video Domains
    - \* News, Surveillance, Traffic, Automobiles, Sports, Gaming, Cooking, HowTo, Travel, Pets, Education, Science, Entertainment, Music, Comedy, Film, Lifestyle, Activism
  - Question Types
    - \* Sequential Understanding: Cooking, How-to, Education
    - \* Predictive Reasoning: Sports, Gaming
    - \* World Knowledge: Science, News
    - \* Causal Reasoning: Surveillance, Activism
    - \* Emotional Reasoning: Entertainment, Film, Comedy
    - \* Analytical Reasoning: Traffic, Automobile

## O.2 Instruction Fine-Tuning Domains

- **Video-Instruct100K**: Video Summarization, Description-based QA (spatial, temporal, relational, reasoning), Creative/Generative QA
- **LLaVA-Video-178K**: Academic Sources, YouTube Videos, ActivityNetQA, NeXT-QA, PerceptionTest, LLaVA-Hound
- **NeXT-QA**: Causal Action Reasoning, Temporal Action Reasoning, Object Interaction Understanding, Daily Activity Scenarios
- **PerceptionTest**: Memory, Abstraction, Physics, Semantics
- **CLEVRER**: Temporal Reasoning, Causal Reasoning, Physical Dynamics, Symbolic Event Representation
- **ActivityNetQA**: Long-Term Spatio-Temporal Reasoning, Complex Web Video Understanding, Human-Annotated QA Pairs

## P Model Size and Budget

We had the following compute budget for our project:

- **Evaluation on ViMUL-Bench (including ablations)**:
  - Total GPU Hours: 9
  - GPU Variant: AMD-MI200 (64GB)
- **Finetuning ViMUL on ViMUL-Instruct**:
  - Total GPU Hours: 200–240
  - GPU Variant: AMD-MI200 (64GB)
- **Translation of 130k samples for cycle consistency using Qwen-3**:
  - Total GPU Hours: 12,000
  - GPU Variant: AMD-MI200 (64GB)

## Q Experimental Setup and Hyperparameters

**Frame Sampling Hyperparameter** In our experimental setup, we investigated the impact of the number of frames used for video representation. We defined  $N$  as the maximum number of frames to be processed per video. While the default setting for ViMUL uses  $N = 32$  uniformly sampled frames at 1 FPS, we also conducted ablation studies with single-frame inputs (e.g., using only the first,

middle, or random frame). Our results showed that such single-frame configurations led to substantial drops in performance, especially in spatio-temporal and reasoning-heavy tasks. Therefore,  $N = 32$  was selected as the optimal configuration, offering a strong balance between computational efficiency and model accuracy.

**Evaluation Hyperparameter** We evaluated the ViMUL-Bench on the following hyperparameters in Tab. 4

Parameter	Value
max_new_tokens	1024
temperature	0
do_sample	False
num_beams	1
batch_size	1

Table 4: Evaluation hyperparameters used during inference.

**Finetuning Hyperparameters** We fine-tune ViMUL-OneVision using the following key hyperparameters:

- Base LLM: DeepSeek-R1-Distill-Llama-8B
- Vision Tower: google/siglip-so400m-patch14-384
- Projector Type: mlp2x\_gelu
- MM Tunable Parts: mm\_mlp\_adapter, mm\_language\_model
- Learning Rate: 1e-5
- Weight Decay: 0.0
- Train Batch Size: 2
- Eval Batch Size: 1
- Gradient Accumulation: 1
- Epochs: 1
- Warmup Ratio: 0.03
- LR Scheduler: Cosine
- Precision: bf16
- Gradient Checkpointing: True
- Max Token Length: 8192

- Torch Compile: True (inductor)
- Deepspeed Config: Zero3
- Save Steps: 1000
- Save Total Limit: 1

## R Instruction to Volunteers

### R.1 Cultural Video Dataset Curation

For each language, we have identified eight distinct cultural categories. We gave the following instructions to the volunteers to collect the videos:

- From your selected language choose 3 to 4 videos from each of the 8 cultural categories
- Please avoid choosing videos that contain sensitive personal information such as private addresses, Social Security numbers, or any other confidential data.
- Whenever possible, select videos filmed in public places to respect privacy. Ensure that the video content does not infringe on the privacy of any individuals or groups.
- Videos should not depict or disclose any private or sensitive content without the consent of the people involved.
- Always ensure that the content of the video respect the privacy of individuals and do not include private or sensitive information.
- If there is any doubt about the appropriateness of a video, please consult the project supervisor or team for clarification.
- Ensure that the video is public and its license is also public.

### R.2 Cultural QA Curation

After the videos are collected, we asked our volunteers to curate QA pairs of both multiple choice and open-ended QA pairs, for each video in their respective native language. The instructions were given as follows:

- Please watch each video carefully. If you feel that the video does not align with the specified category or the native language, you are welcome to replace it with a new YouTube video link. Ensure that the video is public and its license is also public.

- After watching the video, you are required to create 3 Multiple Choice Questions (MCQs) and 1 Short Answer Question (SAQ) in English based on the video's content.
- Additionally, for each question, you must provide a translated version in the native language along with the answer.

### Instructions for Writing Multiple Choice Questions (MCQs)

- **Video-Dependent Questions:** Ensure that each question requires the viewer to watch the video to answer. Avoid questions that could be answered with general knowledge or basic text comprehension. *For example:* "Where is the University of Central Florida?" is incorrect, while "Where is the university shown in the video?" is correct.
- **Contextual Understanding:** The question should assess the viewer's understanding of the video's context, including spatial relationships, object interactions, and scene transitions. *For example:* "Where does the car park after crossing the street, as shown in the video?" is a suitable question.
- **Avoid Ambiguity:** Ensure that all questions and answer options are clear, unambiguous, and directly linked to the video content. Vague or open-ended questions can confuse models and reduce the benchmark's effectiveness. *For example:* Instead of asking "What happens next?", use "What does the person do immediately after sitting down?"
- **Multiple Choices with Plausible Distractors:** Distractors (incorrect options) should be plausible and require watching the video carefully to rule out. Distractors may involve objects or events that appear similar to the correct answer or occur at different times in the video. *For example:* In the question "What does the person eat in the video?", possible choices could be:
  - A. A sandwich
  - B. An apple
  - C. A book (implausible distractor)
  - D. A banana (plausible distractor)
- **Cultural or Contextual Awareness:** Ensure that questions are culturally sensitive and relevant to the video's context, especially if the

video pertains to specific cultural events or actions. This will help avoid bias and make the dataset more generalizable. *For example:* "What festival is being celebrated in the video based on the decorations and activities?"

- **MCQs Format:** The MCQs should be written in the following format:
  - **Question:** <Question>
  - **Answer:** Correct\_Answer (A. Option 1, B. Option 2, C. Option 3, D. Option 4)
- **Example MCQ:**
  - **Question:** What traditional Japanese garment are the individuals in the Video wearing?
  - **Answer:** Kimono (A. Kimono, B. Sari, C. Hanbok, D. Cheongsam)

### R.3 Instructions for Writing Short Answer Questions (SAQs)

- **Concise and Precise Questions:** The questions should be brief and to the point, avoiding unnecessary complexity. *For example:* "What color is the car that drives by at the beginning?"
- **Culturally Relevant Questions:** If the video depicts culturally specific events or actions, ensure that the question is contextually appropriate for that setting. This ensures relevance and avoids bias. *For example:* "What festival is being celebrated in the video?" (based on visual cues like decorations or attire).
- **Answer Length:** Answers should be concise, generally between 1-10 words.

### R.4 Verifying Phi-4 Scores

Volunteers are tasked with verifying the phi-4 scores assigned to model predictions for a set of videos. The verification process involves reviewing the video, the question, the ground truth answer, the model's prediction, and the assigned score. Based on this review, you will determine whether the assigned score is appropriate or if it needs adjustment.

### Step-by-Step Instructions

1. **Watch the Video:** Using the mint\_video\_id, watch the video associated with each QA.

2. **Review the Question, Ground Truth Answer, and Model Prediction:** Carefully read the question (Q), the ground truth answer (A), and the model's prediction (Prediction) in the provided columns.
3. **Check the Assigned Score:** Look at the score already assigned in the Score column. The score ranges from 0 to 5, with increments of 0.5.
4. **Assess the Model's Prediction:** Ask yourself: *Does the model's prediction deserve the score it was given compared to the ground truth answer?*
5. **Select Your Response:**
  - If **YES**, indicating that the model's prediction is appropriately scored, select **YES** from the dropdown menu under the Do you agree with the score? column. Leave the Your Score column empty.
  - If **NO**, indicating that the model's prediction does not deserve the score assigned, select **NO** from the dropdown menu under the Do you agree with the score? column. Then, select your preferred score from the dropdown menu in the Your Score column.

## R.5 Verifying Machine Translated Generic Category QA pairs

Volunteers are tasked with verifying the machine-translated QA pairs for different languages. There are two types of QA pairs: Open-ended and Multiple Choice Questions (MCQs). The goal is to verify the accuracy of the translations for both types of questions and answers. If the translation is correct, the volunteer will mark it as "YES"; if incorrect, they will correct the translation and provide the updated information in the appropriate columns.

### Step-by-Step Instructions

1. **Understand the Columns:** The Excel sheet contains the following columns for English (Ground Truth):
  - English\_Question
  - English\_MCQ\_Choice\_1, English\_MCQ\_Choice\_2, English\_MCQ\_Choice\_3, English\_MCQ\_Choice\_4

- English\_Answer

The following columns are for the translations:

- Translated\_Question
- Translated\_MCQ\_Choice\_1, Translated\_MCQ\_Choice\_2, Translated\_MCQ\_Choice\_3, Translated\_MCQ\_Choice\_4
- Translated\_Answer

2. **Check the Translation Accuracy:** The English text serves as the Ground Truth. For each row, review the translated versions of the question, MCQ choices, and answer.

### 3. Verify the Translation:

- If the translation is **correct**, write YES in the Is the translation correct? column.
- If the translation is **incorrect**, write NO in the Is the translation correct? column, and insert the **correct translation** in the respective columns:
  - Correct\_Translated\_Question
  - Correct\_Translated\_MCQ\_Choice\_1
  - Correct\_Translated\_MCQ\_Choice\_2
  - Correct\_Translated\_MCQ\_Choice\_3
  - Correct\_Translated\_MCQ\_Choice\_4
  - Correct\_Translated\_Answer

### 4. Handling Open-ended QA Pairs:

- For Open-ended QA pairs, the MCQ columns will be empty.
- If the question is Open-ended and the translation is incorrect, only fill in the Correct\_Translated\_Question column and leave the MCQ and answer columns empty.

### Example

#### Incorrect Translation Example:

- **If the translated question is incorrect:**
  - Write NO in the Is the translation correct? column.
  - Insert the correct translation in the Correct\_Translated\_Question column.
  - Leave the other columns (MCQ choices and answer) empty if they are correct.

<b>Languages</b>	<b>Country</b>	<b>Script</b>	<b>Family</b>	<b>Specification</b>
Arabic	UAE, Saudi, Egypt	Arabic	Afro-Asiatic	High
Bengali	Bangladesh, India	Bengali	Indo-European	High
Chinese	China	Chinese	Sino-Tibetan	High
French	France	Latin	Indo-European	High
German	Germany	Latin	Indo-European	High
Hindi	India	Devanagari	Indo-European	High
Japanese	Japan	Kanji/Kana	Japonic	High
Russian	Russia	Cyrillic	Indo-European	High
Sinhala	Sri Lanka	Sinhalese	Indo-European	Low
Spanish	Spain	Latin	Indo-European	High
Swedish	Sweden	Latin	Indo-European	High
Tamil	India	Tamil	Dravidian	Low
Urdu	Pakistan	Arabic	Indo-European	Low

Table 5: Language Classification by Country, Script, Family, and Specification