



# MMVU: MEASURING EXPERT-LEVEL MULTI-DISCIPLINE VIDEO UNDERSTANDING

Yilun Zhao\* Lujing Xie\* Haowei Zhang\* Guo Gan Yitao Long Zhiyuan Hu Tongyan Hu Weiyuan Chen Chuhan Li Junyang Song Zhijian Xu Chengye Wang Weifeng Pan Ziyao Shangguan Xiangru Tang Zhenwen Liang Yixin Liu Chen Zhao Arman Cohan

Yale NLP MMVU Team

## ABSTRACT

We introduce **MMVU**, a comprehensive expert-level, multi-discipline benchmark for evaluating foundation models in video understanding. **MMVU** includes 3,000 expert-annotated questions spanning 27 subjects across four core disciplines: Science, Healthcare, Humanities & Social Sciences, and Engineering. Compared to prior benchmarks, **MMVU** features three key advancements. First, it challenges models to apply domain-specific knowledge and perform expert-level reasoning to analyze specialized-domain videos, moving beyond the basic visual perception typically assessed in current video benchmarks. Second, each example is annotated by human experts from scratch. We implement strict data quality controls to ensure the high quality of the dataset. Finally, each example is enriched with expert-annotated reasoning rationals and relevant domain knowledge, facilitating in-depth analysis. We conduct an extensive evaluation of 32 frontier multimodal foundation models on **MMVU**. The latest System-2-capable models, o1 and Gemini 2.0 Flash Thinking, achieve the highest performance among the tested models. However, they still fall short of matching human expertise. Through in-depth error analyses and case studies, we offer actionable insights for future advancements in expert-level, knowledge-intensive video understanding for specialized domains.



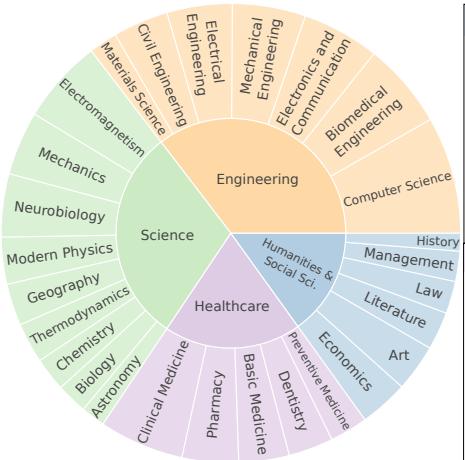
**Project Page:** [mmvu-benchmark.github.io/](https://mmvu-benchmark.github.io/)



**MMVU Data:** [huggingface.co/datasets/yale-nlp/MMVU](https://huggingface.co/datasets/yale-nlp/MMVU)



**MMVU Code:** [github.com/yale-nlp/MMVU](https://github.com/yale-nlp/MMVU)



<b>Science - Chemistry</b>  <p>Q: Assume that 2.24 liters of gas fully participates in the reaction shown in the video under the standard temperature and pressure condition, how many grams of precipitate are produced approximately?</p> <p>(A) 10.0 (B) 5.0 (C) 12.0 (D) 15.0 (E) 20.0</p>	<b>Engineering - Computer Science</b>  <p>Q: Assume that the algorithm is correctly implemented to fully sort the list. What is the total number of comparisons performed in the algorithm?</p> <p><b>Answer: 14</b></p>
<b>Healthcare - Basic Medicine</b>  <p>Q: Which virus does the video depict?</p> <p>(A) Norovirus (B) Measles virus  <b>(C) Hemorrhagic fever virus</b>  (D) Human papillomavirus  (E) Arboviral encephalitis virus</p>	<b>Humanities &amp; Social Science - Art</b>  <p>Q: Which cinematic shooting technique is shown in the video?</p> <p><b>Answer: Dolly Zoom</b></p>

Figure 1: Overview of the **MMVU** benchmark. **MMVU** includes 3,000 expert-annotated examples, covering 27 subjects across four core disciplines. It is specifically designed to assess multimodal foundation models in expert-level, knowledge-intensive video understanding and reasoning tasks.

\*Core Contributors. All authors' contributions are detailed in the Contribution section.

## 1 INTRODUCTION

Foundation models have demonstrated remarkable capabilities in reasoning across various domains, yet their ability to handle expert-level knowledge remains a critical area of evaluation (Hendrycks et al., 2021; Yue et al., 2024a). In recent years, researchers have developed numerous benchmarks to assess these models’ proficiency in specialized domains, primarily focusing on text-based reasoning (Hendrycks et al., 2021; Wang et al., 2024d; Feng et al., 2024; Sun et al., 2024) and image-based contexts (Lu et al., 2024; Yue et al., 2024a;b; Zhang et al., 2024a; Li et al., 2024g). However, as capabilities of foundation models expand across multiple modalities, there is a significant gap in evaluating expert-level reasoning over specialized-domain *videos*. This gap is particularly concerning as video is one of the most information-rich and naturalistic modalities, and is widely used to convey complex, dynamic information in specialized fields like healthcare, engineering, and scientific research (He et al., 2024). Unlike static text or images, expert-level videos often capture temporal dynamics, procedural knowledge, and complex interactions that are essential in many specialized domains. For example, in science, expert-level and knowledge-intensive reasoning might involve analyzing a chemical reaction video (Figure 1). A model must identify key reaction stages based on subtle visual cues like color changes or the formation of precipitates, which requires integrating chemical knowledge in addition to recognizing visual patterns.

To bridge this gap, we introduce **MMVU**, a comprehensive benchmark measuring Multimodal founddation models in expert-level, Multi-discipline Video Understanding and reasoning. **MMVU** consists of 3,000 expert-annotated QA examples over 1,529 specialized-domain videos, spanning 27 subjects across four key disciplines: Science, Healthcare, Humanities & Social Sciences, and Engineering. To ensure both the breadth of domain knowledge and the depth of reasoning required for **MMVU**, we implement a textbook-guided data annotation process. Expert annotators first locate key concepts from textbooks in their fields, then source relevant videos and create corresponding questions that require domain knowledge and expert-level reasoning to comprehend the videos. Each example also includes expert-annotated reasoning rationale and relevant domain knowledge, facilitating fine-grained evaluation of model performance. Thorough data quality controls are implemented to ensure high quality of **MMVU**.

We conduct an extensive evaluation on **MMVU**, covering 32 frontier multimodal foundation models from 17 organizations. Notably, the latest o1 model demonstrates the highest performance among all tested models, approaching the expertise of human experts. Despite this progress, other models still fall noticeably short of human-level capabilities. For instance, GPT-4o achieves a score of 66.7%, which is substantially lower than the benchmark set by human experts (*i.e.*, 86.8%) in the open-book setting. Our analysis highlights the effectiveness of CoT reasoning, which generally enhances model performance compared to directly generating final answers without intermediate reasoning steps. To deepen understanding of the current models’ limitations, we perform an in-depth error analysis of frontier models, including numerous case studies reviewed by human experts. These insights provide valuable guidance for future advancements in the field.

## 2 RELATED WORK

**Video Understanding Benchmark.** Existing video understanding benchmarks primarily focus on *general-purpose* video comprehension tasks, such as action recognition (Heilbron et al., 2015; Sigurdsson et al., 2016; Liu et al., 2020; Deng et al., 2023), captioning and description (Xu et al., 2016; Krishna et al., 2017; Li et al., 2024c; Takahashi et al., 2024; Wu et al., 2021), grounding (Lei et al., 2018; Wang et al., 2022; Chen et al., 2023c; Keser et al., 2023), temporal reasoning (Jang et al., 2017; Liu et al., 2024b; Shangguan et al., 2024; Cores et al., 2024; Cai et al., 2024; Keser et al., 2024; Li et al., 2024e), and long video understanding (Zhang et al., 2023b; Wang et al., 2024b; Nagrani et al., 2024; Atallah et al., 2024; Fang et al., 2024). The rise of video-based foundation models (Tang et al., 2023; Zhang et al., 2023a; Fei et al., 2024; Huang et al., 2024b) has driven the development of new benchmarks that include diverse video comprehension tasks for more comprehensive evaluation (Xiao et al., 2021; Ning et al., 2023; Li et al., 2024d; Fu et al., 2024; Li et al., 2024f; Khattak et al., 2024; Yang et al., 2024b). However, these benchmarks remain predominantly focused on natural scenes and general-purpose tasks. A significant gap persists in benchmarks targeting *expert-level* and *knowledge-intensive reasoning* over specialized-domain videos, where both

Table 1: Comparison between **MMVU** and existing multi-discipline benchmarks for evaluating foundation models. In the “QA Type” column, “MC” denotes Multiple-Choice questions, “Open” denotes Open-ended questions, and “T/F” denotes True-False questions.

<b>Dataset</b>	<b>QA Type</b>	<b>Data Source</b>	<b>College</b>	<b>Detailed Solution</b>	<b>Rational? Knowledge?</b>
			<b>Level?</b>	<b>Level?</b>	
<b><i>Text</i></b>					
MMLU (Hendrycks et al., 2021)	MC	Exam, course, textbook	✓	✗	✗
MMLU-Pro (Wang et al., 2024d)	MC	Datasets → Human & LLM augment	✓	✗	✗
C-Eval (Huang et al., 2023)	MC	Exam	✓	✗	✗
SciEval (Sun et al., 2024)	MC, Open	Internet, datasets → LLM rewrite	✓	✗	✗
TheoremQA (Chen et al., 2023a)	MC, T/F, Open	Internet, exam → Human rewrite	✓	✗	✓
SciKnowEval (Feng et al., 2024)	MC, T/F, Open	Textbooks, database, other datasets → LLM rewrite	✓	✗	✓
<b><i>Text + Image</i></b>					
VisScience (Jiang et al., 2024)	MC, Open	Internet, exam, textbook	✗	✗	✗
EXAMS-V (Das et al., 2024)	MC	Exam	✗	✗	✗
ScienceQA (Lu et al., 2022)	MC	Internet, course	✗	✓	✗
SceMQA (Liang et al., 2024)	MC, Open	Internet, exam	✗	✓	✗
CharXiv (Wang et al., 2024e)	Open	arXiv paper → Human annotate	✓	✗	✗
MMSci (Li et al., 2024g)	MC	Scientific paper → LLM generate	✓	✗	✗
OlympicArena (Huang et al., 2024a)	MC, T/F, Open	Olympic competitions	✓	✓	✗
MMMU (Yue et al., 2024a)	MC, Open	Internet, exam, textbook	✓	17.6%	✗
CMMMU (Zhang et al., 2024a)	MC, T/F, Open	Internet, exam, textbook	✓	2.1%	✗
MMMU-Pro (Yue et al., 2024b)	MC	MMMU → Human & LLM augment	✓	15.4%	✗
<b><i>Text + Video</i></b>					
MMWorld (He et al., 2024)	MC	Human experts (24%) / LLM-gen (76%)	39.5%	✗	✗
<b>MMVU (ours)</b>	MC, Open	Human experts annotate from scratch	✓	✓	✓

visual perception and domain-specific expertise are required—especially in critical fields like healthcare, engineering, and science (He et al., 2024).

**Multi-discipline Evaluation Benchmark.** The rapid development of foundation models has significantly enhanced expert-level reasoning across various disciplines (Touvron et al., 2023; Jiang et al., 2023; Yang et al., 2024a; Google, 2024; OpenAI, 2024b). Early benchmarks focused on domain-specific tasks for textual domains, establishing a foundation for assessing the models’ strengths and limitations in expert reasoning (Welbl et al., 2017; Clark et al., 2018b; Hendrycks et al., 2021; Suzgun et al., 2023; Zhong et al., 2024; Chen et al., 2023a; Wang et al., 2024d; Zhao et al., 2024). More recently, benchmarks have evolved to include multimodal tasks (Yue et al., 2024a; Lu et al., 2024; Zhang et al., 2024a; Yue et al., 2024b; Li et al., 2024g; Wang et al., 2024e), emphasizing visual perception and advanced reasoning with domain knowledge. However, these efforts remain largely limited to *static* images. Developing a high-quality, multidisciplinary video benchmark presents greater challenges than those for text or image-based tasks due to the scarcity of suitable resources (*e.g.*, textbooks or exam questions). This leaves the critical modality of videos and video-based expert-level reasoning significantly underexplored. Recent work, MMWorld (He et al., 2024), has made pioneering strides by incorporating videos across multiple disciplines. However, only a limited portion of its dataset (39.5%) requires domain-specific expertise<sup>1</sup>, and 76.4% of the examples are generated by the GPT-4V model. Moreover, most existing benchmarks provide only the ground-truth answer, restricting researchers’ ability to conduct a fine-grained evaluation. To address this limitation, **MMVU** includes expert-annotated reasoning rationales and relevant domain knowledge for each example, enabling a more nuanced assessment of expert-level reasoning. Table 1 further distinguishes the difference between **MMVU** and existing multi-discipline benchmarks.

### 3 MMVU BENCHMARK

We present **MMVU**, a comprehensive evaluation benchmark that focuses on measuring progress on knowledge-intensive, expert-level reasoning in the video modality. **MMVU** has the following

<sup>1</sup>To estimate the proportion of MMWorld examples requiring domain expertise, we randomly sampled 200 instances from the human-annotated subset and engaged three annotators for evaluation. An example was classified as requiring domain expertise if at least one annotator marked it as such.

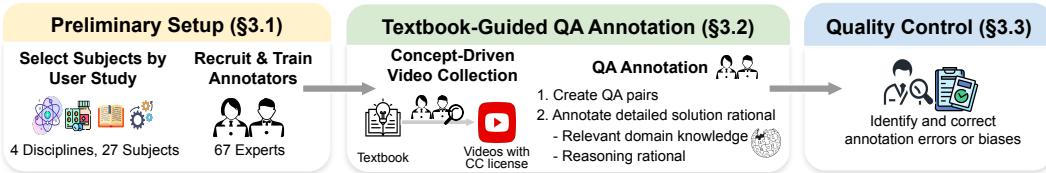


Figure 2: An overview of the **MMVU** benchmark construction pipeline.

key features: (1) **Breadth of Domain Knowledge**: We employ a textbook-guided QA annotation pipeline to ensure the wide coverage of domain knowledge within each subject (§3.2). (2) **Depth of Expert-level Reasoning**: Each example in **MMVU** requires models to comprehend specialized-domain video context, applying expert knowledge and reasoning (§3.2). (3) **True Visual Understanding**: Recent studies (Yue et al., 2024b; Chen et al., 2024a; Zhang et al., 2024b) have shown that visual content is unnecessary for many examples in current multimodal benchmarks. To alleviate this issue, each example in **MMVU** is carefully validated by human experts to confirm that video comprehension is required for accurate answering (§3.3). (4) **Support of Fine-grained Evaluation**: We provide expert-annotated solutions and the requisite knowledge for each example (§3.2), enabling more comprehensive analysis for future research (§4.3). Figure 2 provides an overview of the three stages involved in constructing **MMVU**, which is detailed in the following subsections.

### 3.1 PRELIMINARY SETUP

We first discuss the preliminary setup for data construction.

**Subject Selection.** To ensure a broad and accurate representation of expert-level video understanding across diverse disciplines, we conduct a user study involving 133 college and graduate students for subject selection. We ask them to curate two QA examples requiring expert-level video understanding in subjects relevant to their field of study, and provide feedback on their experiences during the curation process. Such a user study-guided approach helps us identify subjects within each discipline that may not be obvious from a top-down selection process. It also offers insights into the challenges of designing expert-level video examples, helping us design and refine the textbook-guided QA annotation process (detailed in §3.2). The authors manually analyze the collected examples and select **27 subjects** (as listed in Figure 1) across four disciplines that align best with our benchmark’s construction desiderata discussed earlier.

**Expert Annotator Recruitment and Training.** For each subject, we assign at least two annotators with relevant expertise. We include a total of 67 expert annotators (detailed biographies are presented in Appendix A.1), comprising 22 third- or fourth-year undergraduate students, 36 graduate students, and nine of the authors. All the annotators also participated in our initial user study. Each annotator is required to finish a training session to learn the annotation protocol (detailed in Appendix A.3) before official annotation.

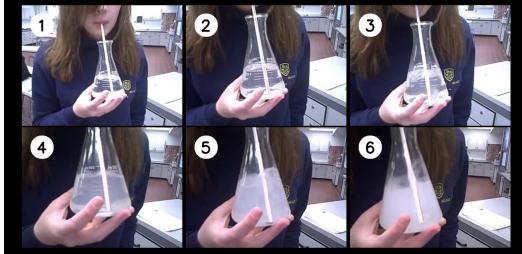
### 3.2 TEXTBOOK-GUIDED QA EXAMPLE ANNOTATION

Constructing a high-quality, expert-level, multi-disciplinary benchmark for video-based tasks is more challenging than the ones for text- or image-based, as there is no existing resources (*e.g.*, textbooks or exam questions) that can be adapted from and each example has to be curated from scratch. Therefore, it is crucial to establish a structured approach that ensures the quality and comprehensiveness of the benchmark. We employ a textbook-guided example annotation pipeline designed to capture both the *breadth of knowledge* and *depth of reasoning*. In brief, annotators first identify key concepts from the textbook and locate relevant videos that align with these concepts. The textbooks for each subject (listed in Appendix A.2) are selected by expert annotators and are recognized as authoritative references in their respective fields. Annotators then curate QA examples and detailed solution rationales. We detail the annotation procedure as follows:

**Concept-Driven CC-Licensed Video Collection.** Annotators are instructed to first review each chapter of the textbook to identify key concepts that inherently require dynamic visual representa-

**Question:**  
Assume that 2.24 liters of gas fully participates in the reaction shown in the video under the standard temperature and pressure condition, how many grams of precipitate are produced approximately?

**Options:** (A) 10.0 ✓ (B) 5.0 (C) 12.0 (D) 15.0 (E) 20.0



**Textbook used for annotation:** "Chemistry, 2nd Edition (Paul Flowers, Klaus Theopold, Richard Langley, William R. Robinson)"

**Annotated Relevant Domain Knowledge (Wikipedia page):**

1. Calcium hydroxide: [https://en.wikipedia.org/wiki/Calcium\\_hydroxide](https://en.wikipedia.org/wiki/Calcium_hydroxide)  
"...When carbon dioxide is passed through limewater, the solution takes on a milky appearance due to precipitation of insoluble calcium carbonate:  $\text{Ca(OH)}_2(\text{aq}) + \text{CO}_2(\text{g}) \rightarrow \text{CaCO}_3(\text{s}) + \text{H}_2\text{O}(\text{l})$ ..."
2. Carbon dioxide: [https://en.wikipedia.org/wiki/Carbon\\_dioxide](https://en.wikipedia.org/wiki/Carbon_dioxide)
3. Ideal gas law: [https://en.wikipedia.org/wiki/Ideal\\_gas\\_law](https://en.wikipedia.org/wiki/Ideal_gas_law)

**Annotated Reasoning Rationale:**  
In the video, a person exhales gas that is continuously introduced into a clear solution, gradually forming a white precipitate. This indicates that the substances involved in the reaction are  $\text{CO}_2$  and limewater. The chemical reaction equation is:  $\text{Ca(OH)}_2 + \text{CO}_2 \rightarrow \text{CaCO}_3 + \text{H}_2\text{O}$ . At the STP, 2.24 liters of  $\text{CO}_2$  corresponds to 0.1 Moles. From balanced equation, 0.1 moles of  $\text{CO}_2$  produce 0.1 moles of  $\text{CaCO}_3$ . Given  $\text{Ca} = 40 \text{ g/mol}$ ,  $\text{C} = 12 \text{ g/mol}$ ,  $\text{O} = 16 \text{ g/mol}$ , the molar mass of  $\text{CaCO}_3 = 40 + 12 + 16 * 3 = 100 \text{ g/mol}$ . Therefore, the mass of  $\text{CaCO}_3 = 0.1 * 100 = 10 \text{ g}$ .

Figure 3: A dataset example from **MMVU** with the discipline of chemistry. Each example in **MMVU** includes expert annotation of relevant domain knowledge and step-by-step reasoning rational.

tion, such as experimental procedures in science or mechanical operations in engineering. They then search for related videos on YouTube having Creative Commons license<sup>2</sup> that effectively illustrate the selected concept. To ensure the collected videos effectively challenge the model’s visual reasoning capabilities, the video should be vision-intensive, requiring models to focus solely on visual information for comprehension. To this end, we ensure that audio tracks are excluded to eliminate potential shortcuts models might exploit through auditory cues; and the video should contain minimal on-screen text, as an overabundance of text may detract from the core visual understanding task. Consequently, videos such as lecture recordings, which typically include slides or text-based explanations that simplify the task of answering associated questions, are excluded.

**QA Annotation.** After identifying suitable videos, annotators are required to create two or three questions, either multiple-choice or open-ended. Each question is designed to test the model’s expert-level reasoning by applying domain-specific knowledge to interpret the video content and derive a solution. Annotators are also required to specify the start and end timestamps of the video clip relevant to answering each question. For annotating multi-choice question, the annotators are required to carefully craft the four distractor options to reflect common misconceptions or plausible alternatives, ensuring that models cannot easily eliminate incorrect options without reasoning over video content. Once the five options are finalized, the annotation interface randomly shuffles them.

**Solution Rationale Annotation.** For each annotated question, annotators must also provide detailed solution for the correct answers. As shown in Figure 3, the solution comprises two key components: (1) *relevant domain knowledge*, which includes a list of domain-specific concepts or keywords necessary for answering the question, with each concept linked to its corresponding Wikipedia page. (2) *reasoning rationale*, which details the step-by-step reasoning process to reach the correct answer. These solution annotations are critical for enhancing transparency in the evaluation process and facilitating future research focused on understanding model failure modes.

### 3.3 DATA QUALITY CONTROL

We next discuss our methods to ensure high data quality.

**Time-Based Annotation Compensation.** As discussed earlier, annotating examples for **MMVU** can be particularly time-intensive, especially when there is limited availability of videos with Creative Commons licenses in the required subjects. To accommodate this and ensure a high-quality

<sup>2</sup>The Creative Commons license enables reusers to distribute, remix, adapt, and build upon the material in any medium or format, so long as attribution is given to the creator. We use YouTube Data API v3 (<https://developers.google.com/youtube/v3>) to verify the license type. Existing video benchmarks typically utilize YouTube videos, yet do not confine their selections to content with CC licenses, introducing potential copyright concerns. We recognize that by restricting our selection to CC-licensed content, we are compelled to forgo coverage of certain subjects (e.g., sports), where CC-licensed videos is scarce.

Table 2: Key statistics of the **MMVU** benchmark.

Statistics	Value
Total Questions	3,000
Validation Set	1,000
Test Set	2,000
Unique Videos	1,529
Video Length (Seconds, avg/max)	51.4 / 228
Number of Disciplines	4
Number of Subjects	27
Multiple Choice Questions	1,858
Question Length (avg/max)	16.8 / 70
Single Choice Length (avg/max)	7.6 / 42
Number of Choices per Question	5
Open-ended Questions	1,142
Question Length (avg/max)	16.4 / 39
Ground-truth Answer Length (avg/max)	1.5 / 7
Number of Required Knowledge per Question (avg/max)	4.3 / 7
Solution Rationale Length (avg/max)	56.6 / 193
Total Number of Unique Knowledge ( <i>i.e.</i> , Wikipedia pages)	4,770

benchmark, we compensate annotators based on the time they spend rather than the number of examples completed, preventing them from rushing through tasks (See Appendix A.5 for annotation compensation details). On average, annotating one example takes 20 minutes and 17 seconds, while validation requires 4 minutes and 12 seconds.

**Human Expert Validation.** To ensure that the final dataset remains high-quality and meets expert-level standards without introducing unnecessary biases, each example in **MMVU** undergoes expert review by one of the authors or top-performing annotators to verify the accuracy of its annotations. Recent studies (Yue et al., 2024b; Chen et al., 2024a; Zhang et al., 2024b; Shangguan et al., 2024) have shown that visual content is unnecessary for many examples in current multimodal benchmarks. To address this concern, each example in **MMVU** is carefully validated by human experts to ensure that video comprehension is required for accurate answering. If an example is determined to be answerable solely through the textual components of the question, a single video frame, or if it contains annotation errors, evaluators first attempt to revise the example. If revision is not feasible, detailed feedback is provided to the original annotator, who then revises and submits it for a second iteration. A total of 523 examples were revised during the data validation process. Among them, 72 examples were still found to be misaligned with our design criteria and were excluded from the final benchmark. Overall,  $1 - \frac{523}{3,000+72} = 83.0\%$  of the initial examples met our design criteria without requiring revisions, indicating the high quality of initial annotation.

### 3.4 MMVU BENCHMARK ANALYSIS

**Data Statistics.** Table 2 presents the key statistics of **MMVU**. It consists of 3,000 examples, which are randomly divided into two subsets: validation and test. The validation set contains 1,000 examples, and is intended for model development and validation. The test set, comprising the remaining 2,000 examples, is strictly reserved for standard evaluation to prevent data contamination (Jacovi et al., 2023; Deng et al., 2024; Glazer et al., 2024). To further promote fair benchmarking, the test set remains hidden. We are developing an online evaluation pipeline on a public platform, enabling researchers to benchmark their models and participate in a public leaderboard.

**Human Performance.** To provide a rough but informative estimate of human-level performance on **MMVU**, we randomly sampled 30 questions per discipline from the test set, resulting in a total of 120 questions for evaluation. Five participants—three graduate students specializing in biology, anesthesiology, and East-Asian literature, along with two of the authors—individually answered these questions. The evaluation proceeded in three phases: (1) **Closed-book Setting:** In the first

phase, participants had 3.5 hours to answer questions without access to external resources. The average accuracy across the four participants was 49.7%. (2) **Open-book Setting**: In the second phase, participants were permitted to use external resources (*e.g.*, internet and textbooks) to review answers they felt uncertain about. They were not informed of the correctness of their initial responses, and a 4-hour time limit was set. This open-book approach led to an increase in average accuracy to 86.8%. (3) **Oracle Setting**: Finally, participants were required to revise each incorrect answer based on ground-truth domain knowledge and self-sourced online resources. The average accuracy after this final revision was 95.3%.

## 4 EXPERIMENTS

This section discusses the experiment setup and our key findings.

### 4.1 EXPERIMENT SETUP

**Evaluated Multimodal Foundation Models.** To establish a comprehensive understanding of the challenges posed by **MMVU** and provide reference points for future research, we evaluate a broad range of frontier multimodal foundation models that support *video* or *multiple images* as input. Specifically, we evaluate **16 series of open-source models**, including InternVL-2 & 2.5 (Chen et al., 2023b; 2024b), Qwen2-VL (Wang et al., 2024a; Yang et al., 2024a), LLaVA-NeXT (Liu et al., 2024a), Pixtral (MistralAI, 2024), DeepSeek-VL2 (Wu et al., 2024), H2OVL Mississippi (Galib et al., 2024), Idefics2 (Laurençon et al., 2024), Aria (Li et al., 2025), LLaVA-NeXT-Video (Li et al., 2024b), LLaVA-OneVision (Li et al., 2024a), Llama-3.2-Vision (Dubey et al., 2024), Phi-3.5-Vision (Abdin et al., 2024), InternVideo2 (Wang et al., 2024c), and VideoLLaMA2 & 2.1 (Cheng et al., 2024). We also evaluate **eight series of proprietary models**, including OpenAI o1 (OpenAI, 2024a) and GPT-4o (OpenAI, 2024b), Gemini-1.5 & 2 and Gemini-Thinking (Google, 2024), GLM-4V-Plus (GLM et al., 2024), Grok-2-Vision (xAI, 2024), and Claude-3.5 (Anthropic, 2024). For open-source models, we prioritize the vLLM pipeline (Kwon et al., 2023) for model inference; otherwise, we use the Transformers pipeline (Wolf et al., 2020). We use the official API service for proprietary models. For models without native video support, following VideoMME (Fu et al., 2024), we provide visual input using the maximum number of images that fits within the model’s context window. §B.1 details the parameter settings and model configurations. We evaluate the models with both **Direct Answer** and **Chain-of-Thought** prompts (presented in appendix B.2), which is adapted from the versions used in MMMU-Pro (Yue et al., 2024b).

**Accuracy Evaluation.** We use accuracy as the primary metric to evaluate model performance on **MMVU**. Following recent benchmarks for foundation model evaluation (Wang et al., 2024e; Lu et al., 2024; He et al., 2024), we employ GPT-4o to assess accuracy. Specifically, given a question, its ground truth answer, and the model’s response, GPT-4o is instructed to extract the final answer from the model response and determine its correctness. The evaluation prompts for both multiple-choice and open-ended questions are presented in Appendix B.3.

### 4.2 MAIN FINDINGS

Table 3 presents the evaluated models’ CoT performance on the **MMVU** benchmark, while Figure 4 illustrates a comparison between the model performance in CoT reasoning and direct answering. Our key findings are as follows:

**MMVU presents substantial challenges for current multimodal foundation models.** Even the top-performing model falls well short of human expert performance. For instance, GPT-4o achieves 66.7% accuracy with CoT prompting, significantly lower than the 86.8% accuracy achieved by human experts in an open-book setting. Notably, while GPT-4o has narrowed the performance gap with human experts in text-based expert-level reasoning on MMLU (88.7% vs 89.8% (Hendrycks et al., 2021)) and image-based expert-level reasoning on MMMU (69.1% vs 82.6% (Yue et al., 2024a)), the gap remains large on **MMVU**. This disparity underscores **MMVU**’s critical role in advancing and evaluating multimodal foundation models’ capabilities in video-based expert reasoning across specialized domains.

Table 3: Accuracy of evaluated foundation models on the **MMVU** validation and test sets using CoT prompts. Model performance is ranked based on overall results on the test set. \*: For o1, as the API access for its multimodal version has not been granted, we randomly sampled 100 examples from the validation set and 200 examples (50 for each core discipline) from the test set. The model’s performance was manually evaluated on Jan 10, 2025, using CoT prompts on ChatGPT platform.

Release	Test Set				Avg. Validation	Avg. Test
	Science	Healthcare	Human. & Social Sci.	Engineering		
<b>Human Performance</b>						
Human Oracle	95.3	93.3	96.0	96.7	95.3	
Human Open-book	86.7	84.7	92.7	83.3	86.8	
Human Closed-book	54.7	42.7	44.7	56.7	49.7	
<b>Proprietary Models</b>						
o1*	2024-12	80.0	78.0	76.0	74.0	79.0
Gemini 2.0 Flash Thinking	2024-12	69.3	71.2	73.4	67.3	69.1
GPT-4o	2024-08	67.2	71.8	72.0	61.6	67.4
Gemini 2.0 Flash	2024-12	70.8	62.7	71.6	63.0	65.9
Gemini 1.5 Pro	2024-09	67.2	68.1	67.0	62.8	65.4
Claude 3.5 Sonnet	2024-10	60.5	64.0	70.9	64.5	65.2
Grok-2-Vision	2024-12	60.6	72.5	72.0	57.4	62.7
GPT-4o-mini	2024-07	60.3	60.9	70.6	59.3	61.6
Gemini 1.5 Flash	2024-09	56.8	57.3	66.3	58.2	58.8
GLM-4V-Plus	2025-01	52.2	57.3	64.9	55.4	56.2
<b>Open-sourced Models</b>						
Qwen2-VL-72B	2024-09	48.0	53.6	61.7	53.9	53.2
DeepSeek-VL2	2024-12	50.3	53.4	58.9	48.6	52.1
InternVL2.5-38B	2024-11	50.3	45.6	52.8	52.8	50.5
Aria	2024-11	46.8	43.3	61.0	49.9	49.3
Llama-3.2-90B-Vision	2024-09	46.5	43.5	53.9	48.1	47.1
DeepSeek-VL2-Small	2024-12	47.5	48.7	47.5	45.1	46.9
Qwen2-VL-7B-Instruct	2024-08	43.6	42.5	43.6	41.2	42.1
InternVL2.5-8B	2024-11	39.2	36.8	47.2	42.3	41.1
VideoLLaMA2.1-7B	2024-10	35.3	38.9	45.4	41.6	39.5
Llama-3.2-11B-Vision	2024-09	40.5	39.4	44.0	35.7	38.9
Phi-3.5-Vision	2024-08	38.3	29.5	45.4	41.1	38.1
LLaVA-OneVision-7B	2024-09	34.3	38.6	40.8	38.8	37.9
Qwen2-VL-2B	2024-08	32.6	40.9	40.4	35.7	36.5
InternVL2-8B	2024-06	36.7	32.9	36.9	37.2	36.3
Idefics3-8B	2024-08	37.0	35.5	44.0	31.2	35.3
VideoLLaMA2-7B	2024-06	32.3	27.7	44.3	35.7	34.4
DeepSeek-VL2-Tiny	2024-12	34.3	33.4	35.8	30.1	33.0
Pixtral-12B	2024-09	36.1	24.6	37.9	30.8	32.3
LLaVA-NeXT-Video-34B	2024-06	31.8	24.6	35.8	30.3	30.5
InternVideo2-8B	2024-08	29.6	31.1	37.2	26.5	29.9
H2OVL Mississippi-2B	2024-10	29.1	29.5	29.4	28.0	29.1
LLaVA-NeXT-Video-7B	2024-06	27.0	31.1	27.3	29.5	28.6

**Performance of open-sourced models.** As for open-source multimodal foundation models, they still lag behind the proprietary models. However, the Qwen2-VL-72B and DeepSeek-VL2 models have achieved performance levels that exceed human benchmarks in closed-book settings and are approaching the performance of leading proprietary models. These advancements highlight the significant progress being made in the development of open-source models.

**CoT reasoning generally improves model performance compared to directly outputting the answer.** However, the degree of improvement varies across different foundation models. For instance, Claude 3.5 Sonnet demonstrated a remarkable enhancement, achieving a notable performance gain of 11.0%, as corroborated by the findings in MMMU-Pro (Yue et al., 2024b).

Conversely, models like GPT-4o exhibited only marginal improvements. These results indicate that the impact of CoT reasoning is not uniformly beneficial across all models on **MMVU**.

**System-2 thinking demonstrates effectiveness.** Models capable of System-2 thinking and employing long CoT demonstrate significant performance advantages. Notably, the o1 and Gemini 2.0 Flash Thinking models achieved the top two results on **MMVU**, illustrating that increasing test-time compute and applying long CoT can significantly enhance model performance in expert-level video reasoning tasks. These results highlight the potential of developing open-source models designed to facilitate and advance System-2 thinking capabilities.

#### 4.3 QUALITATIVE ANALYSIS

To gain a deeper understanding of the capabilities and limitations of frontier models on **MMVU**, we perform comprehensive case studies and error analysis by humans. The inclusion of expert-annotated reasoning rationales and domain knowledge for each example in **MMVU** facilitate a more effective analysis compared to datasets that provide only answers. We focus on four top-performing models, GPT-4o, Qwen2-VL-72B, Llama-3.2-90B-Vision, and DeepSeek-VL2, for human evaluation. From the **MMVU** validation set, we randomly sample 50 error cases for each model. These cases are analyzed by the authors using ground-truth features (*i.e.*, expert-annotated reasoning rationales and required domain knowledge) as references. We identify following six primary errors:

**Visual Perception Error (18%):** The model fails to accurately interpret spatial, temporal, or semantic aspects of visual information within a video. Additionally, it might “hallucinate”, detecting objects or events that are not actually present in the video. Figure 5 (left) is a typical related instance where the model fails to correctly perceive the traversal order of the binary tree. Similarly, Figure 18 shows that the model mistakenly identifies the device shell in the video as water, leading to completely wrong reasoning about the device’s function.

**Misuse or Lack Domain Knowledge in Visual Perception (20%):** The model fails to apply the domain-specific expertise required to accurately interpret specialized concepts or elements within the video. For example, in a medical video, it may identify objects but fail to recognize their technical terms or misunderstand their importance within the procedure being demonstrated. Moreover, as shown in Figure 20, the model correctly perceives the ascending numbers (array indices), but misuses its pretrained knowledge and misidentifies them as the numbers to be sorted. It leads to the wrong conclusion that the video demonstrates a sorting algorithm. This limitation underscores a gap in the model’s ability to integrate domain knowledge with visual perception effectively.

**Misuse or Lack Domain Knowledge in Reasoning (27%):** The model fails to effectively recall and apply domain knowledge during its reasoning processes. For instance, when addressing questions over chemistry videos, it may fail to correctly apply relevant chemical equations, leading to errors in computing the reaction mass. A notable example is Figure 5 (right), where the model misuses the domain knowledge that bats often live in unsanitary environments and makes the wrong inference that poor hygiene conditions are the cause of virus outbreaks. Besides, in Figure 25, the model lacks the domain knowledge about relevant chemical equations, so that it cannot correctly answer the question. This limitation underscores the model’s inability to integrate domain knowledge into its reasoning processes effectively.

**Heavy Reliance on Textual Information (20%):** The model predominantly depends on textual information for problem-solving, especially when addressing multiple-choice questions, as it evaluates each option individually without leveraging the actual video content. For instance, Figure 26 shows the model ignores the video information about the reason of the disease and overly focuses on the textual question. Similar limitations have been observed in other multimodal benchmarks (Fu

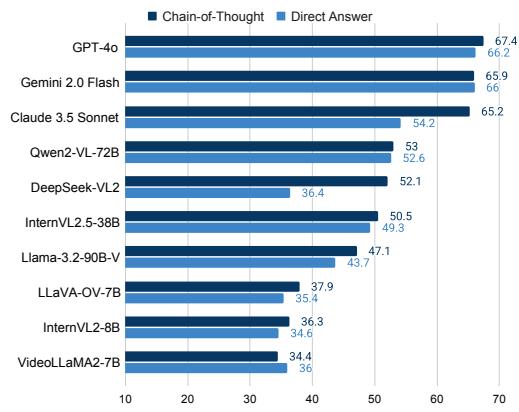


Figure 4: Comparison of model performance between CoT and direct answering on the validation set. The full results are provided in §C.1.

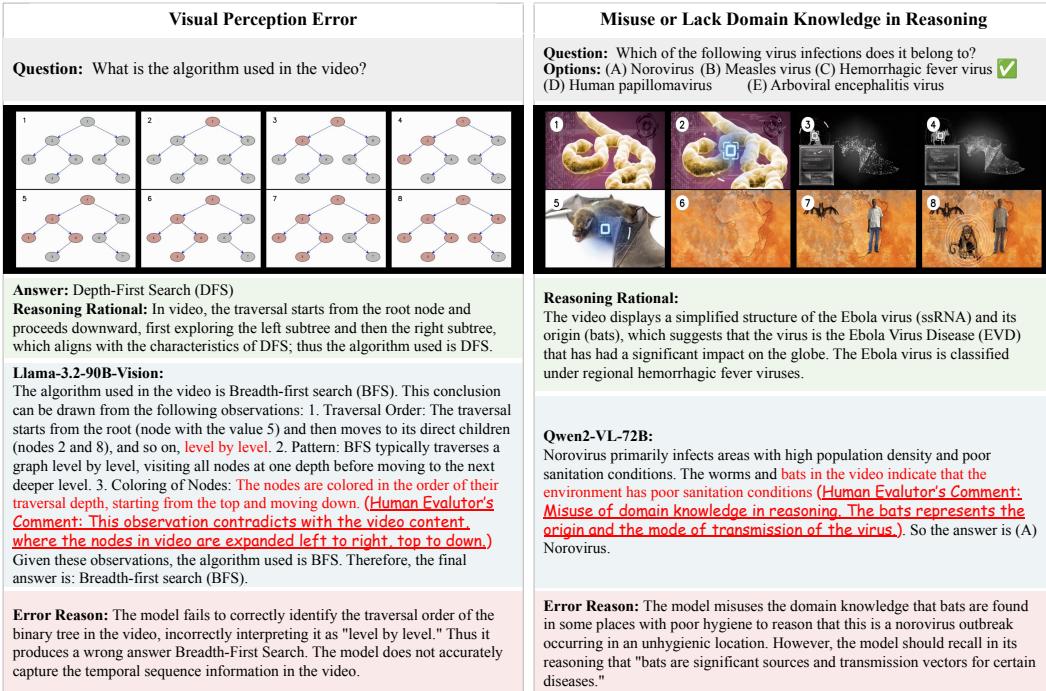


Figure 5: Illustrations of visual perception error and misuse or lack domain knowledge in reasoning.

et al., 2024; Yue et al., 2024a). This gap suggests future work in enhancing multimodal reasoning by more effectively incorporating non-textual content into the reasoning process.

**Logical Reasoning Error (6%):** The model exhibits inconsistencies between its reasoning process and final answer, leading to self-contradiction. As depicted in Figure 28, the analysis of one specific option contradicts with the other reasoning steps, which is a typical self-contradiction logical error.

**Other Error (9%):** This includes other errors, such as refusing to answer a question due to insufficient context or safety concerns, generating a response that exceeds the output limit, generating repetitive information, or making incorrect math computation.

## 5 CONCLUSION

We introduce **MMVU**, a high-quality, multi-disciplinary benchmark designed to assess the expert-level, knowledge-intensive reasoning capabilities of multimodal foundation models on specialized-domain videos. Each example in **MMVU** is annotated by human experts from scratch. We employ a textbook-guided example annotation pipeline designed to capture both the breadth of knowledge and depth of reasoning. In our evaluation of 32 frontier multimodal foundation models, we find that while the latest o1 model achieves the highest performance among all tested models—approaching human expert-level proficiency—a notable performance gap remains between other models and human experts. Additionally, models employing CoT reasoning consistently outperform those that generate final answers directly. Through comprehensive error analysis and case studies, we identify persistent challenges of **MMVU**, offering valuable insights for advancing foundation models' capabilities to achieve expert-level video understanding in specialized domains.

## AUTHOR CONTRIBUTION

The author contributions are summarized below:

- **Project Lead:** Yilun Zhao
- **Project Conception:** Yilun Zhao, Lujing Xie, Yitao Long, Zhiyuan Hu, Zhenwen Liang, Xiangru Tang, Yixin Liu, Chen Zhao, Arman Cohan
- **User Study:** Every author
- **Data Annotation Protocol Development:** Yilun Zhao, Lujing Xie, Chengye Wang
- **Data Annotation Task Management:** Lujing Xie, Haowei Zhang
- **Data Annotation:** Lujing Xie, Haowei Zhang, Tongyan Hu, Weiyuan Chen, Junyang Song, Zhijian Xu, Weifeng Pan, Guo Gan, Yitao Long
- **Data Validation:** Lujing Xie, Haowei Zhang, Tongyan Hu, Weiyuan Chen, Yilun Zhao, Junyang Song
- **Data Annotation Expense:** Yilun Zhao
- **Codebases and Results:** Yilun Zhao, Guo Gan
- **Error Analysis and Case Study:** Haowei Zhang, Lujing Xie, Yilun Zhao, Weiyuan Chen
- **Manuscript Writing:** Yilun Zhao, Haowei Zhang, Arman Cohan
- **Manuscript Editing:** Every author

## REFERENCES

Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dong-dong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi, Xin Jin, Nikos Karampatziakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong Luo, Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro Mendes, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo de Rosa, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Jilong Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan, Chenruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. Phi-3 technical report: A highly capable language model locally on your phone, 2024. URL <https://arxiv.org/abs/2404.14219>.

Bruce Alberts, Alexander Johnson, Julian Lewis, Martin Raff, Keith Roberts, and Peter Walter. *Molecular Biology of the Cell*. Garland Science, 6th edition, 2014.

Phillip E Allen and Douglas R Holberg. *CMOS analog circuit design*. Elsevier, 2011.

Anthropic. Introducing the next generation of claude, 2024. URL <https://www.anthropic.com/news/claude-3-family>.

Mumtaz Anwar, Riyaz Ahmad Rather, and Zeenat Farooq. *Fundamentals and advances in medical biotechnology*. Springer, 2022.

- Steven Ascher and Edward Pincus. *The Filmmaker's Handbook: A Comprehensive Guide for the Digital Age*. Plume, Penguin Random House, 5th edition, 2012.
- Kirolos Atallah, Chenhui Gou, Eslam Abdelrahman, Khushbu Pahwa, Jian Ding, and Mohamed Elhoseiny. Infinibench: A comprehensive benchmark for large multimodal models in very long video understanding, 2024. URL <https://arxiv.org/abs/2406.19875>.
- Peter William Atkins, Julio De Paula, and James Keeler. *Atkins' physical chemistry*. Oxford university press, 2023.
- Eugene A. Avallone, Theodore Baumeister, and Ali M. Sadegh. *Marks' Standard Handbook for Mechanical Engineers*. McGraw-Hill Education, 12th edition, 2018.
- Ashwani Bedi and Ramsey Dabby. *Structure for Architects: A Case Study in Steel, Wood, and Reinforced Concrete Design*. Routledge, 1st edition, 2019.
- Fred G Bell. *Engineering geology and construction*. CRC Press, 2004.
- Olivier Blanchard. *Macroeconomics*. Pearson, 9th edition, 2024.
- David S. Bright, Anastasia H. Cortes, et al. *Principles of Management*. OpenStax, Rice University, 2019. Available at <https://openstax.org/details/books/principles-management>.
- Theodore L. Brown, H. Eugene LeMay, Bruce E. Bursten, Catherine J. Murphy, Patrick M. Woodward, and Matthew E. Stoltzfus. *Chemistry: The Central Science*. Pearson, 15th edition, 2023.
- Laurence L. Brunton, Randa Hilal-Dandan, and Bjorn Knollman. *Goodman & Gilman's: The Pharmacological Basis of Therapeutics*. McGraw-Hill Education, 13th edition, 2017.
- Randal E Bryant and David Richard O'Hallaron. *Computer systems: a programmer's perspective*. Prentice Hall, 2011.
- Mu Cai, Reuben Tan, Jianrui Zhang, Bocheng Zou, Kai Zhang, Feng Yao, Fangrui Zhu, Jing Gu, Yiwu Zhong, Yuzhang Shang, Yao Dou, Jaden Park, Jianfeng Gao, Yong Jae Lee, and Jianwei Yang. Temporalbench: Benchmarking fine-grained temporal understanding for multimodal video models, 2024. URL <https://arxiv.org/abs/2410.10818>.
- William D Callister Jr and David G Rethwisch. *Materials science and engineering: an introduction*. John wiley & sons, 2020.
- Krishan K. Chawla. *Composite Materials: Science and Engineering*. Springer, 3rd edition, 2012.
- Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, and Feng Zhao. Are we on the right way for evaluating large vision-language models?, 2024a. URL <https://arxiv.org/abs/2403.20330>.
- Wenhu Chen, Ming Yin, Max Ku, Pan Lu, Yixin Wan, Xueguang Ma, Jianyu Xu, Xinyi Wang, and Tony Xia. TheoremQA: A theorem-driven question answering dataset. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 7889–7901, Singapore, December 2023a. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.489. URL <https://aclanthology.org/2023.emnlp-main.489>.
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *arXiv preprint arXiv:2312.14238*, 2023b.
- Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, Ji Ma, Jiaqi Wang, Xiaoyi Dong, Hang Yan, Hewei Guo, Conghui He, Botian Shi, Zhenjiang Jin, Chao Xu, Bin Wang, Xingjian Wei, Wei Li, Wenjian Zhang, Bo Zhang, Pinlong Cai, Licheng Wen, Xiangchao Yan, Min Dou, Lewei Lu, Xizhou Zhu, Tong Lu, Dahua Lin, Yu Qiao, Jifeng Dai, and Wenhai Wang. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites, 2024b. URL <https://arxiv.org/abs/2404.16821>.

- Zhihong Chen, Ruifei Zhang, Yibing Song, Xiang Wan, and Guanbin Li. Advancing visual grounding with scene knowledge: Benchmark and method. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 15039–15049, June 2023c.
- Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, and Lidong Bing. Videollama 2: Advancing spatial-temporal modeling and audio understanding in video-lmms. *arXiv preprint arXiv:2406.07476*, 2024. URL <https://arxiv.org/abs/2406.07476>.
- Mary Ann Clark, Jung Choi, and Matthew Douglas. *Biology*. OpenStax, Rice University, 2nd edition, 2018a. Available at <https://openstax.org/details/books/biology-2e>.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018b.
- Jonathan Clayden, Nick Greeves, and Stuart Warren. *Organic chemistry*. Oxford University Press, USA, 2012.
- Daniel Cores, Michael Dorkenwald, Manuel Mucientes, Cees G. M. Snoek, and Yuki M. Asano. Tybench: Redesigning video-language evaluation, 2024. URL <https://arxiv.org/abs/2410.07752>.
- Thomas H Cormen, Charles E Leiserson, Ronald L Rivest, and Clifford Stein. *Introduction to algorithms*. MIT press, 2022.
- Braja M. Das. *Principles of Geotechnical Engineering*. Cengage Learning, 9th edition, 2017.
- Rocktim Jyoti Das, Simeon Emilov Hristov, Haonan Li, Dimitar Iliyanov Dimitrov, Ivan Koychev, and Preslav Nakov. Exams-v: A multi-discipline multilingual multimodal exam benchmark for evaluating vision language models, 2024. URL <https://arxiv.org/abs/2403.10378>.
- Mackenzie L. Davis and David A. Cornwell. *Introduction to Environmental Engineering*. McGraw-Hill Education, 5th edition, 2012.
- Andong Deng, Taojiannan Yang, and Chen Chen. A large-scale study of spatiotemporal representation learning with a new benchmark on action recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 20519–20531, October 2023.
- Chunyuan Deng, Yilun Zhao, Yuzhao Heng, Yitong Li, Jiannan Cao, Xiangru Tang, and Arman Cohan. Unveiling the spectrum of data contamination in language model: A survey from detection to remediation. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics: ACL 2024*, pp. 16078–16092, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.951. URL <https://aclanthology.org/2024.findings-acl.951>.
- Avi Domb, Boaz Mizrahi, and Shady Farah. *Biomaterials and Biopolymers*. Springer, 2023.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, More, and Zhiwei Zhao. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- John D. Enderle and Joseph D. Bronzino. *Introduction to Biomedical Engineering*. Academic Press, 4th edition, 2017.
- Xinyu Fang, Kangrui Mao, Haodong Duan, Xiangyu Zhao, Yining Li, Dahua Lin, and Kai Chen. Mmbench-video: A long-form multi-shot benchmark for holistic video understanding, 2024. URL <https://arxiv.org/abs/2406.14515>.
- Adam Feather, David Randall, and Mona Waterhouse. *Kumar and Clark's Clinical Medicine E-Book: Kumar and Clark's Clinical Medicine E-Book*. Elsevier Health Sciences, 2020.

- Jiajun Fei, Dian Li, Zhidong Deng, Zekun Wang, Gang Liu, and Hui Wang. Video-ccam: Enhancing video-language understanding with causal cross-attention masks for short and long videos, 2024. URL <https://arxiv.org/abs/2408.14023>.
- Kehua Feng, Keyan Ding, Weijie Wang, Xiang Zhuang, Zeyuan Wang, Ming Qin, Yu Zhao, Jianhua Yao, Qiang Zhang, and Huajun Chen. Sciknoweval: Evaluating multi-level scientific knowledge of large language models, 2024. URL <https://arxiv.org/abs/2406.09098>.
- Harry L Field and John M Long. *Introduction to agricultural engineering technology: a problem solving approach*. Springer, 2018.
- Paul Flowers, Klaus Theopold, Richard Langley, and William R. Robinson. *Chemistry*. OpenStax, Rice University, 2nd edition, 2019. Available at <https://openstax.org/details/books/chemistry-2e>.
- Erin H Fouberg and Alexander B Murphy. *Human Geography: People, Place, and Culture*. John Wiley & Sons, 2020.
- Fabrizio Frigeni. *Industrial Robotics Control: Mathematical Models, Software Architecture, and Electronics Design*. Springer, 2022.
- Victoria Fromkin, Robert Rodman, and Nina Hyams. *An Introduction to Language*. Cengage Learning, 11th edition, 2017.
- Chaoyou Fu, Yuhua Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, Peixian Chen, Yanwei Li, Shaohui Lin, Sirui Zhao, Ke Li, Tong Xu, Xiawu Zheng, Enhong Chen, Rongrong Ji, and Xing Sun. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis, 2024. URL <https://arxiv.org/abs/2405.21075>.
- Shaikat Galib, Shanshan Wang, Guanshuo Xu, Pascal Pfeiffer, Ryan Chesler, Mark Landry, and Sri Satish Ambati. H2ovl-mississippi vision language models technical report, 2024. URL <https://arxiv.org/abs/2410.13611>.
- Elliot Glazer, Ege Erdil, Tamay Besiroglu, Diego Chicharro, Evan Chen, Alex Gunning, Caroline Falkman Olsson, Jean-Stanislas Denain, Anson Ho, Emily de Oliveira Santos, Olli Järvinieniemi, Matthew Barnett, Robert Sandler, Matej Vrzala, Jaime Sevilla, Qiuyu Ren, Elizabeth Pratt, Lionel Levine, Grant Barkley, Natalie Stewart, Bogdan Grechuk, Tetiana Grechuk, Shreepranav Varma Enugandla, and Mark Wildon. Frontiermath: A benchmark for evaluating advanced mathematical reasoning in ai, 2024. URL <https://arxiv.org/abs/2411.04872>.
- Team GLM, :, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Dan Zhang, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui, Jie Tang, Jing Zhang, Jingyu Sun, Juanzi Li, Lei Zhao, Lindong Wu, Lucen Zhong, Mingdao Liu, Minlie Huang, Peng Zhang, Qinkai Zheng, Rui Lu, Shuaiqi Duan, Shudan Zhang, Shulin Cao, Shuxun Yang, Weng Lam Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, Xiaohan Zhang, Xiaotao Gu, Xin Lv, Xinghan Liu, Xinyi Liu, Xinyue Yang, Xixuan Song, Xunkai Zhang, Yifan An, Yifan Xu, Yilin Niu, Yuantao Yang, Yueyan Li, Yushi Bai, Yuxiao Dong, Zehan Qi, Zhaoyu Wang, Zhen Yang, Zhengxiao Du, Zhenyu Hou, and Zihan Wang. Chatglm: A family of large language models from glm-130b to glm-4 all tools, 2024. URL <https://arxiv.org/abs/2406.12793>.
- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- Google. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context, 2024.
- Nikolai V. Gorbunov. *Tissue Barriers in Disease, Injury and Regeneration*. Elsevier, 1st edition, 2022.
- Steven A. Greenlaw, David Shapiro, and Daniel MacDonald. *Principles of Economics*. OpenStax, Rice University, 3rd edition, 2023. Available at <https://openstax.org/details/books/principles-economics-3e>.

- David J Griffiths. *Introduction to electrodynamics*. Cambridge University Press, 2023.
- Allan R Hambley. *Electrical Engineering: Principles and Applications*. Pearson London, UK, 2018.
- Xuehai He, Weixi Feng, Kaizhi Zheng, Yujie Lu, Wanrong Zhu, Jiachen Li, Yue Fan, Jianfeng Wang, Linjie Li, Zhengyuan Yang, Kevin Lin, William Yang Wang, Lijuan Wang, and Xin Eric Wang. Mmworld: Towards multi-discipline multi-faceted world model evaluation in videos, 2024. URL <https://arxiv.org/abs/2406.08407>.
- Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 961–970, 2015. doi: 10.1109/CVPR.2015.7298698.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=d7KBjmI3GmQ>.
- Darrel Hess and Tom L. McKnight. *McKnight's Physical Geography: A Landscape Appreciation*. Pearson, 13th edition, 2021.
- HLTCOE@JHU. Turkle: A web-based tool for managing annotation tasks. <https://github.com/hltcoe/turkle>, 2024. Accessed: 2024-11-01.
- Paul Horowitz and Winfield Hill. The art of electronics, 2015.
- Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Jiayi Lei, Yao Fu, Maosong Sun, and Junxian He. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models, 2023. URL <https://arxiv.org/abs/2305.08322>.
- Zhen Huang, Zengzhi Wang, Shijie Xia, Xuefeng Li, Haoyang Zou, Ruijie Xu, Run-Ze Fan, Lyumanshan Ye, Ethan Chern, Yixin Ye, Yikai Zhang, Yuqing Yang, Ting Wu, Binjie Wang, Shichao Sun, Yang Xiao, Yiyuan Li, Fan Zhou, Steffi Chern, Yiwei Qin, Yan Ma, Jiadi Su, Yixiu Liu, Yuxiang Zheng, Shaoting Zhang, Dahua Lin, Yu Qiao, and Pengfei Liu. Olympicarena: Benchmarking multi-discipline cognitive reasoning for superintelligent ai, 2024a. URL <https://arxiv.org/abs/2406.12753>.
- Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, Yaohui Wang, Xinyuan Chen, Limin Wang, Dahua Lin, Yu Qiao, and Ziwei Liu. Vbench: Comprehensive benchmark suite for video generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 21807–21818, June 2024b.
- Peter Huber and Alastair Mullis. *The CISG: A new textbook for students and practitioners*. Sellier de Gruyter, 2009.
- Alon Jacovi, Avi Caciularu, Omer Goldman, and Yoav Goldberg. Stop uploading test data in plain text: Practical strategies for mitigating data contamination by evaluation benchmarks. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 5075–5084, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.308. URL <https://aclanthology.org/2023.emnlp-main.308/>.
- Yunseok Jang, Yale Song, Youngjae Yu, Youngjin Kim, and Gunhee Kim. Tgif-qa: Toward spatio-temporal reasoning in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2758–2766, 2017.

Albert Qiaochu Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L'elio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b. *ArXiv*, abs/2310.06825, 2023. URL <https://api.semanticscholar.org/CorpusID:263830494>.

Zhihuan Jiang, Zhen Yang, Jinhao Chen, Zhengxiao Du, Weihan Wang, Bin Xu, Yuxiao Dong, and Jie Tang. Visscience: An extensive benchmark for evaluating k12 educational multi-modal scientific reasoning, 2024. URL <https://arxiv.org/abs/2409.13730>.

Eric R. Kandel, James H. Schwartz, Thomas M. Jessell, Steven A. Siegelbaum, and A.J. Hudspeth. *Principles of Neural Science*. McGraw-Hill Education, 6th edition, 2021.

Ilker Keser, Andrea Pedrotti, Mustafa Dogan, Michele Cafagna, Emre Can Acikgoz, Letitia Parcalabescu, Iacer Calixto, Anette Frank, Albert Gatt, Aykut Erdem, and Erkut Erdem. Vilma: A zero-shot benchmark for linguistic and temporal grounding in video-language models, 2023. URL <https://arxiv.org/abs/2311.07022>.

Ilker Keser, Andrea Pedrotti, Mustafa Dogan, Michele Cafagna, Emre Can Acikgoz, Letitia Parcalabescu, Iacer Calixto, Anette Frank, Albert Gatt, Aykut Erdem, and Erkut Erdem. ViLMA: A zero-shot benchmark for linguistic and temporal grounding in video-language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=liuqDwmbQJ>.

Muhammad Uzair Khattak, Muhammad Ferjad Naeem, Jameel Hassan, Muzammal Naseer, Federico Tombari, Fahad Shahbaz Khan, and Salman Khan. How good is my video lmm? complex video reasoning and robustness evaluation suite for video-lmms, 2024. URL <https://arxiv.org/abs/2405.03690>.

Richard R. Kibbe, Roland O. Meyer, John E. Neely, and Warran T. White. *Machine Tool Practices*. Pearson, 11th edition, 2019.

David R. Klein. *Organic Chemistry as a Second Language: First Semester Topics*. John Wiley & Sons, 2024.

Fred S. Kleiner. *Art Through the Ages: A Global History, Volume I*. Cengage Learning, 16th edition, 2020.

Ann Kordas, Ryan J. Lynch, et al. *World History Volume 1*. OpenStax, Rice University, 2022. Available at <https://openstax.org/details/books/world-history-volume-1>.

Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning events in videos. In *Proceedings of the IEEE international conference on computer vision*, pp. 706–715, 2017.

Vinay Kumar, Abul K. Abbas, and Jon C. Aster. *Robbins and Cotran Pathologic Basis of Disease*. Elsevier, 10th edition, 2020.

Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.

Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. What matters when building vision-language models?, 2024. URL <https://arxiv.org/abs/2405.02246>.

Jie Lei, Licheng Yu, Mohit Bansal, and Tamara Berg. TVQA: Localized, compositional video question answering. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun'ichi Tsujii (eds.), *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 1369–1379, Brussels, Belgium, October-November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1167. URL <https://aclanthology.org/D18-1167/>.

- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer, 2024a. URL <https://arxiv.org/abs/2408.03326>.
- Dongxu Li, Yudong Liu, Haoning Wu, Yue Wang, Zhiqi Shen, Bowen Qu, Xinyao Niu, Fan Zhou, Chengan Huang, Yanpeng Li, Chongyan Zhu, Xiaoyi Ren, Chao Li, Yifan Ye, Peng Liu, Lihuan Zhang, Hanshu Yan, Guoyin Wang, Bei Chen, and Junnan Li. Aria: An open multimodal native mixture-of-experts model, 2025. URL <https://arxiv.org/abs/2410.05993>.
- Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *arXiv preprint arXiv:2407.07895*, 2024b.
- Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, Limin Wang, and Yu Qiao. Mvbench: A comprehensive multi-modal video understanding benchmark, 2024c. URL <https://arxiv.org/abs/2311.17005>.
- Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, et al. Mvbench: A comprehensive multi-modal video understanding benchmark. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22195–22206, 2024d.
- Shicheng Li, Lei Li, Shuhuai Ren, Yuanxin Liu, Yi Liu, Rundong Gao, Xu Sun, and Lu Hou. Vitatecs: A diagnostic dataset for temporal concept understanding of video-language models, 2024e. URL <https://arxiv.org/abs/2311.17404>.
- Xinhao Li, Zhenpeng Huang, Jing Wang, Kunchang Li, and Limin Wang. Videoeval: Comprehensive benchmark suite for low-cost evaluation of video foundation model, 2024f. URL <https://arxiv.org/abs/2407.06491>.
- Zekun Li, Xianjun Yang, Kyuri Choi, Wanrong Zhu, Ryan Hsieh, HyeonJung Kim, Jin Hyuk Lim, Sungyoung Ji, Byungju Lee, Xifeng Yan, Linda Ruth Petzold, Stephen D. Wilson, Woosang Lim, and William Yang Wang. Mmsci: A dataset for graduate-level multi-discipline multimodal scientific understanding, 2024g. URL <https://arxiv.org/abs/2407.04903>.
- Zhenwen Liang, Kehan Guo, Gang Liu, Taicheng Guo, Yujun Zhou, Tianyu Yang, Jiajun Jiao, Renjie Pi, Jipeng Zhang, and Xiangliang Zhang. Scemqa: A scientific college entrance level multimodal question answering benchmark, 2024. URL <https://arxiv.org/abs/2402.05138>.
- Samuel J. Ling, Jeff Sanny, and William Moebs. *University Physics Volume 1*. OpenStax, Rice University, 2016a. Available at <https://openstax.org/details/books/university-physics-volume-1>.
- Samuel J. Ling, Jeff Sanny, and William Moebs. *University Physics Volume 2*. OpenStax, Rice University, 2016b. Available at <https://openstax.org/details/books/university-physics-volume-2>.
- Samuel J. Ling, Jeff Sanny, and William Moebs. *University Physics Volume 3*. OpenStax, Rice University, 2016c. Available at <https://openstax.org/details/books/university-physics-volume-3>.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024a. URL <https://llava-v1.github.io/blog/2024-01-30-llava-next/>.
- Jiaying Liu, Sijie Song, Chunhui Liu, Yanghao Li, and Yueyu Hu. A benchmark dataset and comparison study for multi-modal human action analytics. *ACM Trans. Multimedia Comput. Commun. Appl.*, 16(2), May 2020. ISSN 1551-6857. doi: 10.1145/3365212. URL <https://doi.org/10.1145/3365212>.
- Yuanxin Liu, Shicheng Li, Yi Liu, Yuxiang Wang, Shuhuai Ren, Lei Li, Sishuo Chen, Xu Sun, and Lu Hou. TempCompass: Do video LLMs really understand videos? In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics: ACL 2024*, pp. 8731–8772, Bangkok, Thailand, August 2024b. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.517. URL <https://aclanthology.org/2024.findings-acl.517/>.

William Lowrie and Andreas Fichtner. *Fundamentals of geophysics*. Cambridge university press, 2020.

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

Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=KUNzEQMWU7>.

Liquin Luo. *Principles of neurobiology*. Garland Science, 2020.

Marina MacKay. *The Cambridge introduction to the novel*. Cambridge University Press, 2010.

Upamanyu Madhow. *Introduction to communication systems*. Cambridge University Press, 2014.

PK Mallick. Fiber-reinforced composites: Materials, manufacturing, and design, 2007.

Gregory N. Mankiw. *Principles of Microeconomics*. Cengage Learning, 9th edition, 2020.

Kenneth Fuller Maxcy, Milton Joseph Rosenau, John M Last, Robert B Wallace, Neal Kohatsu, and Ross Brownson. *Maxcy-Rosenau-Last public health & preventive medicine*. McGraw-Hill, 2008.

MistralAI. Announcing pixtral 12b, 2024. URL <https://mistral.ai/news/pixtral-12b/>.

Arsha Nagrani, Mingda Zhang, Ramin Mehran, Rachel Hornung, Nitesh Bharadwaj Gundavarapu, Nilpa Jha, Austin Myers, Xingyi Zhou, Boqing Gong, Cordelia Schmid, Mikhail Sirotenko, Yukun Zhu, and Tobias Weyand. Neptune: The long orbit to benchmarking long video understanding. 2024.

Jill Nelmes (ed.). *Introduction to Film Studies*. Routledge, 5th edition, 2012.

Donald A Nield and Adrian Bejan. *Convection in Porous Media*. Springer, 2017.

Munan Ning, Bin Zhu, Yujia Xie, Bin Lin, Jiaxi Cui, Lu Yuan, Dongdong Chen, and Li Yuan. Video-bench: A comprehensive benchmark and toolkit for evaluating video-based large language models, 2023. URL <https://arxiv.org/abs/2311.16103>.

Katsuhiko Ogata. *Modern Control Engineering*. Prentice Hall, 5th edition, 2010.

OpenAI. Openai o1 system card. 2024a. URL <https://api.semanticscholar.org/CorpusID:274611667>.

OpenAI. Hello gpt-4o, 2024b. URL <https://openai.com/index/hello-gpt-4o/>.

Judith A. Owen, Jenni Punt, and Sharon A. Stranford. *Kuby Immunology*. W.H. Freeman, 8th edition, 2018.

David A. Patterson and John L. Hennessy. *Computer organization and design: The hardware/software interface*. Elsevier, 6th edition, 2022.

Onno Rudolf Pols. *Stellar structure and evolution*. Astronomical Institute Utrecht NY, 2011.

Dale Purves, GJ Augustine, David Fitzpatrick, WC Hall, AS LaMantia, RD Mooney, ML Platt, and LE White. Neuroscience (sixth edit), 2018.

C Gonzalez Rafael and E Woods Richard. *Digital Image Processing*. Pearson Education, 2018.

- Colin Renfrew and Paul Bahn. *Archaeology: Theories, Methods, and Practice*. Thames & Hudson, 7th edition, 2016.
- Robert E. Ricklefs. *The Economy of Nature*. W.H. Freeman, 7th edition, 2013.
- Barbara Ryden and Bradley M Peterson. *Foundations of astrophysics*. Cambridge University Press, 2020.
- Daniel V. Schroeder. *An introduction to thermal physics*. Oxford University Press, 2020.
- Robert Sedgewick and Kevin Wayne. Algorithms (4th edn). *Google Scholar Google Scholar Digital Library Digital Library*, 2011.
- Ziyao Shangguan, Chuhan Li, Yuxuan Ding, Yanan Zheng, Yilun Zhao, Tesca Fitzgerald, and Arman Cohan. Tomato: Assessing visual temporal reasoning capabilities in multimodal foundation models, 2024. URL <https://arxiv.org/abs/2410.23266>.
- Gunnar A Sigurdsson, Güll Varol, Xiaolong Wang, Ali Farhadi, Ivan Laptev, and Abhinav Gupta. Hollywood in homes: Crowdsourcing data collection for activity understanding. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I* 14, pp. 510–526. Springer, 2016.
- Abraham Silberschatz, Peter B. Galvin, and Greg Gagne. *Operating System Concepts*. John Wiley & Sons, 10th edition, 2018.
- Liangtai Sun, Yang Han, Zihan Zhao, Da Ma, Zhennan Shen, Baocai Chen, Lu Chen, and Kai Yu. Scieval: A multi-level large language model evaluation benchmark for scientific research. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(17):19053–19061, Mar. 2024. doi: 10.1609/aaai.v38i17.29872. URL <https://ojs.aaai.org/index.php/AAAI/article/view/29872>.
- Mirac Suzgun, Nathan Scales, Nathanael Schärlí, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Akanksha Chowdhery, Quoc Le, Ed Chi, Denny Zhou, and Jason Wei. Challenging BIG-bench tasks and whether chain-of-thought can solve them. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 13003–13051, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.824. URL <https://aclanthology.org/2023.findings-acl.824/>.
- Rikito Takahashi, Hirokazu Kiyomaru, Chenhui Chu, and Sadao Kurohashi. Abstractive multi-video captioning: Benchmark dataset construction and extensive evaluation. In Nicoletta Calzolari, Min-Yen Kan, Veronique Hoste, Alessandro Lenci, Sakriani Sakti, and Nianwen Xue (eds.), *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pp. 57–69, Torino, Italia, May 2024. ELRA and ICCL. URL <https://aclanthology.org/2024.lrec-main.5>.
- Yunlong Tang, Jing Bi, Siting Xu, Luchuan Song, Susan Liang, Teng Wang, Daoan Zhang, Jie An, Jingyang Lin, Rongyi Zhu, Ali Vosoughi, Chao Huang, Zeliang Zhang, Feng Zheng, Jianguo Zhang, Ping Luo, Jiebo Luo, and Chenliang Xu. Video understanding with large language models: A survey. *arXiv preprint arXiv:2312.17432*, 2023.
- Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mi-haylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Bin Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez,

- Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. *ArXiv*, abs/2307.09288, 2023. URL <https://api.semanticscholar.org/CorpusID:259950998>.
- Chris Turner. *Contract law*. Routledge, 2013.
- Ray Turner. *Arbitration awards: a practical approach*. John Wiley & Sons, 2008.
- G Cornelis Van Kooten. *Land resource economics and sustainable development: economic policies and the common good*. UBC Press, 2011.
- Hal R. Varian. *Intermediate Microeconomics: A Modern Approach*. W.W. Norton & Company, 8th edition, 2010.
- William R Wagner, Shelly E Sakiyama-Elbert, Guigen Zhang, and Michael J Yaszemski. *Biomaterials Science: An Introduction to Materials in Medicine*. Elsevier, 2020.
- Jinfeng Wang. *Intelligent Manufacturing System and Intelligent Workshop*. Springer.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution, 2024a. URL <https://arxiv.org/abs/2409.12191>.
- Weihan Wang, Zehai He, Wenyi Hong, Yean Cheng, Xiaohan Zhang, Ji Qi, Shiyu Huang, Bin Xu, Yuxiao Dong, Ming Ding, and Jie Tang. Lvbench: An extreme long video understanding benchmark, 2024b. URL <https://arxiv.org/abs/2406.08035>.
- Yi Wang, Kunchang Li, Xinhao Li, Jiashuo Yu, Yinan He, Chenting Wang, Guo Chen, Baoqi Pei, Ziang Yan, Rongkun Zheng, Jilan Xu, Zun Wang, Yansong Shi, Tianxiang Jiang, Songze Li, Hongjie Zhang, Yifei Huang, Yu Qiao, Yali Wang, and Limin Wang. Internvideo2: Scaling foundation models for multimodal video understanding, 2024c. URL <https://arxiv.org/abs/2403.15377>.
- Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhua Chen. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark, 2024d. URL <https://arxiv.org/abs/2406.01574>.
- Yuxuan Wang, Difei Gao, Licheng Yu, Weixian Lei, Matt Feiszli, and Mike Zheng Shou. Geb+: A benchmark for generic event boundary captioning, grounding and retrieval. In *European Conference on Computer Vision*, pp. 709–725. Springer, 2022.
- Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi Wu, Haotian Liu, Sadhika Malladi, Alexis Chevalier, Sanjeev Arora, and Danqi Chen. Charxiv: Charting gaps in realistic chart understanding in multimodal llms, 2024e. URL <https://arxiv.org/abs/2406.18521>.
- Johannes Welbl, Nelson F. Liu, and Matt Gardner. Crowdsourcing multiple choice science questions. In Leon Derczynski, Wei Xu, Alan Ritter, and Tim Baldwin (eds.), *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pp. 94–106, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/W17-4413. URL <https://aclanthology.org/W17-4413>.
- Edward J. Wing and Fred J. Schiffman. *Cecil Essentials of Medicine*. Elsevier, 10th edition, 2021.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierrick Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 38–45, Online, October 2020. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/2020.emnlp-demos.6>.

- Bo Wu, Shoubin Yu, Zhenfang Chen, Josh Tenenbaum, and Chuang Gan. Star: A benchmark for situated reasoning in real-world videos. In J. Vanschoren and S. Yeung (eds.), *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, volume 1, 2021. URL [https://datasets-benchmarks-proceedings.neurips.cc/paper\\_files/paper/2021/file/5ef059938ba799aaa845e1c2e8a762bd-Paper-round2.pdf](https://datasets-benchmarks-proceedings.neurips.cc/paper_files/paper/2021/file/5ef059938ba799aaa845e1c2e8a762bd-Paper-round2.pdf).
- Zhiyu Wu, Xiaokang Chen, Zizheng Pan, Xingchao Liu, Wen Liu, Damai Dai, Huazuo Gao, Yiyang Ma, Chengyue Wu, Bingxuan Wang, Zhenda Xie, Yu Wu, Kai Hu, Jiawei Wang, Yaofeng Sun, Yukun Li, Yishi Piao, Kang Guan, Aixin Liu, Xin Xie, Yuxiang You, Kai Dong, Xingkai Yu, Haowei Zhang, Liang Zhao, Yisong Wang, and Chong Ruan. Deepseek-vl2: Mixture-of-experts vision-language models for advanced multimodal understanding, 2024. URL <https://arxiv.org/abs/2412.10302>.
- xAI. Grok-2 beta release, 2024. URL <https://x.ai/blog/grok-2>.
- Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. Next-qa: Next phase of question-answering to explaining temporal actions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9777–9786, June 2021.
- Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5288–5296, 2016.
- John A Yagiela, Frank J Dowd, Bart Johnson, Angelo Mariotti, and Enid A Neidle. *Pharmacology and Therapeutics for Dentistry-E-Book: Pharmacology and Therapeutics for Dentistry-E-Book*. Elsevier Health Sciences, 2010.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. Qwen2 technical report, 2024a. URL <https://arxiv.org/abs/2407.10671>.
- Jihan Yang, Shusheng Yang, Anjali W. Gupta, Rilyn Han, Li Fei-Fei, and Saining Xie. Thinking in space: How multimodal large language models see, remember, and recall spaces, 2024b. URL <https://arxiv.org/abs/2412.14171>.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen. 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 (CVPR)*, pp. 9556–9567, June 2024a. URL [https://openaccess.thecvf.com/content/CVPR2024/html/Yue\\_MMMU\\_A\\_Massive\\_Multi-discipline\\_Multimodal\\_Understanding\\_and\\_Reasoning\\_Benchmark\\_for\\_CVPR\\_2024\\_paper.html](https://openaccess.thecvf.com/content/CVPR2024/html/Yue_MMMU_A_Massive_Multi-discipline_Multimodal_Understanding_and_Reasoning_Benchmark_for_CVPR_2024_paper.html).
- Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang, Kai Zhang, Shengbang Tong, Yuxuan Sun, Botao Yu, Ge Zhang, Huan Sun, Yu Su, Wenhui Chen, and Graham Neubig. Mmmu-pro: A more robust multi-discipline multimodal understanding benchmark, 2024b. URL <https://arxiv.org/abs/2409.02813>.
- Ge Zhang, Xinrun Du, Bei Chen, Yiming Liang, Tongxu Luo, Tianyu Zheng, Kang Zhu, Yuyang Cheng, Chunpu Xu, Shuyue Guo, Haoran Zhang, Xingwei Qu, Junjie Wang, Ruibin Yuan, Yizhi Li, Zekun Wang, Yudong Liu, Yu-Hsuan Tsai, Fengji Zhang, Chenghua Lin, Wenhao Huang, and Jie Fu. Cmmmu: A chinese massive multi-discipline multimodal understanding benchmark, 2024a. URL <https://arxiv.org/abs/2401.11944>.

- Hang Zhang, Xin Li, and Lidong Bing. Video-LLaMA: An instruction-tuned audio-visual language model for video understanding. In Yansong Feng and Els Lefever (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 543–553, Singapore, December 2023a. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-demo.49. URL <https://aclanthology.org/2023.emnlp-demo.49/>.
- Hongjie Zhang, Yi Liu, Lu Dong, Yifei Huang, Zhen-Hua Ling, Yali Wang, Limin Wang, and Yu Qiao. Movqa: A benchmark of versatile question-answering for long-form movie understanding, 2023b. URL <https://arxiv.org/abs/2312.04817>.
- Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Peng Gao, and Hongsheng Li. Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems?, 2024b. URL <https://arxiv.org/abs/2403.14624>.
- Yilun Zhao, Hongjun Liu, Yitao Long, Rui Zhang, Chen Zhao, and Arman Cohan. Financemath: Knowledge-intensive math reasoning in finance domains. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 12841–12858, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.693. URL <https://aclanthology.org/2024.acl-long.693/>.
- Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. AGIEval: A human-centric benchmark for evaluating foundation models. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Findings of the Association for Computational Linguistics: NAACL 2024*, pp. 2299–2314, Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-naacl.149. URL <https://aclanthology.org/2024.findings-naacl.149/>.

## Appendix Contents

<b>A MMVU Preliminary Setup</b>	<b>24</b>
A.1 Annotator Biography . . . . .	24
A.2 Textbook for Each Subject . . . . .	26
A.3 Annotation Guideline and Interface . . . . .	29
A.4 Validation Guideline and Interface . . . . .	31
A.5 Data Annotation and Validation Payment . . . . .	31
<b>B Experiment Setup</b>	<b>32</b>
B.1 Configuration of Evaluated Models . . . . .	32
B.2 Chain-of-Thought and Direct Answer Prompts . . . . .	33
B.3 Prompts for Accuracy Evaluation . . . . .	35
<b>C Experiment</b>	<b>36</b>
C.1 Comparison Between CoT Reasoning and Direct Answering . . . . .	36
C.2 Error Case Analysis: Visual Perception Error . . . . .	37
C.3 Error Case Analysis: Misuse or Lack Domain Knowledge in Visual Perception . . . . .	40
C.4 Error Case Analysis: Misuse or Lack Domain Knowledge in Reasoning . . . . .	43
C.5 Error Case Analysis: Heavy Reliance on Textual Information . . . . .	46
C.6 Error Case Analysis: Logical Reasoning Error . . . . .	48

## A MMVU PRELIMINARY SETUP

### A.1 ANNOTATOR BIOGRAPHY

The detailed biographies of the annotators involved in **MMVU** construction are presented in Table 4. All annotators are from universities ranked in the Top 500 of the 2024 QS Global Rankings<sup>3</sup> and are fluent in English.

ID	Year	Major	Assigned Subject(s)	Author?	Validator?
1	1st year Master	Biomedical Engineering	Biomedical Engineering Computer Science Electrical Engineering	X	X
2	1st year Master	Bioinformatics	Biomedical Engineering	X	X
3	1st year Master	Biological Engineering	Biomedical Engineering	X	X
4	2nd year Master	Biomedical Engineering	Biomedical Engineering Electronics and Communication	X	X
5	5th year PhD	Agricultural and Biosystems Engineering	Biomedical Engineering	X	X
6	2nd year Master	Architecture	Civil Engineering	X	X
7	3rd year PhD	Civil Engineering	Civil Engineering Mechanical Engineering	X	X
8	–	–	–	✓	✓
9	3rd year Undergraduate	Electrical Engineering	Computer Science Electrical Engineering	X	X
10	2nd year Master	Electrical Engineering	Computer Science Electronics and Communication	X	X
11	2nd year Master	Electrical Engineering	Computer Science Mechanical Engineering	X	X
12	3rd year Undergraduate	Software Engineering	Computer Science	X	X
13	2nd year Master	Computer Science	Computer Science	X	X
14	–	–	– Electrical Engineering Computer Science Electronics and Communication	✓	X
15	1st year PhD	Electrical Engineering	Electrical Engineering	X	X
16	1st year PhD	Electrical Engineering	Electronics and Communication	X	X
17	–	–	– Electrical Engineering	✓	✓
18	1st year Master	Electrical Engineering	Electrical Engineering Mechanical Engineering	X	X
19	1st year PhD	Electrical Engineering	Electronics and Communication	X	X
20	3rd year PhD	Food Science	Mechanics	X	X
21	4th year PhD	Materials Science	Materials Science	X	X
22	4th year Undergraduate	Aerospace Engineering	Materials Science Mechanical Engineering	X	X
23	4th year Undergraduate	Mechanical Engineering	Materials Science Mechanical Engineering	X	✓
24	2nd year PhD	Mechanical Engineering	Mechanical Engineering	X	X
25	1st year PhD	Mechanical Engineering	Mechanical Engineering	X	X
26	1st year Master	Medicine	Basic Medicine Clinical Medicine	X	X
27	1st year Master	Radiology	Basic Medicine Clinical Medicine	X	X
28	1st year Master	Dentistry	Basic Medicine Dentistry	X	X
29	1st year PhD	Nursing	Basic Medicine Pharmacy	X	X
30	3rd year Undergraduate	Epidemiology	Basic Medicine Preventive Medicine	X	X
31	3rd year Undergraduate	Medicine	Clinical Medicine	X	X
32	–	–	– Clinical Medicine Pharmacy	✓	✓
33	2nd year PhD	Medicine	Clinical Medicine Pharmacy	X	X

Table 4: Biographies of 73 annotators involved in **MMVU** construction (Author biographies are hidden to protect identity confidentiality).

<sup>3</sup><https://www.topuniversities.com/world-university-rankings>

ID	Year	Major	Assigned Subject(s)	Author?	Validator?
34	4th year PhD	Dentistry	Dentistry	X	X
35	3rd year Undergraduate	Dentistry	Dentistry	X	X
36	4th year PhD	Dentistry	Dentistry	X	X
37	1st year PhD	Public Health	Pharmacy Preventive Medicine	X	X
38	4th year Undergraduate	Pharmacy	Pharmacy	X	X
39	3rd year PhD	East Asian Studies	Art	X	X
40	4th year PhD	Literature	Art History Literature	X	X
41	-	-	- History	✓	X
42	1st year PhD	Economics	Economics	X	X
43	4th year Undergraduate	Accounting	Economics Law	X	X
44	4th year PhD	Finance	Economics	X	X
45	3rd year PhD	Public Administration	Law Management	X	X
46	1st year Master	Literature	Literature	X	X
47	5th year PhD	Linguistics	Literature	X	X
48	3rd year Undergraduate	Public Administration	Management	X	X
49	5th year PhD	Astronomy	Astronomy	X	X
50	-	-	-	✓	✓
51	2nd year Master	Astronomy	Astronomy	X	X
52	-	-	- Geography	✓	X
53	3rd year PhD	Biology	Biology	X	X
54	1st year PhD	Biology	Biology Neurobiology	X	X
55	3rd year PhD	Marine Biology	Biology Chemistry	X	X
56	-	-	- Chemistry	✓	X
57	1st year PhD	Chemistry	Chemistry	X	X
58	3rd year Undergraduate	Chemistry	Chemistry	X	X
59	1st year PhD	Physics	Electromagnetism	X	X
60	4th year Undergraduate	Physics	Electromagnetism Thermodynamics	X	X
61	4th year PhD	Physics	Electromagnetism	X	X
62	1st year PhD	Physics	Electromagnetism Mechanics	X	X
			Thermodynamics		
63	1st year Master	Physics	Thermodynamics	X	X
64	3rd year Undergraduate	Agricultural and Environmental Sciences	Electromagnetism		
65	4th year PhD	Physics	Geography Thermodynamics	X	X
			Mechanics		
			Modern Physics		
66	1st year PhD	Physics	Mechanics	X	X
67	3rd year PhD	Physics	Mechanics	X	X
68	4th year PhD	Physics	Modern Physics	X	X
69	3rd year Undergraduate	Neurobiology	Neurobiology	X	X
70	1st year PhD	Neurobiology	Neurobiology	X	X
71	-	-	- Neurobiology	✓	✓
72	3rd year Undergraduate	Biology	Neurobiology	X	X
73	1st year Master	Biology	Neurobiology	X	X

Table 5: Biographies of 73 annotators involved in MMVU construction (Author biographies are hidden to protect identity confidentiality).

## A.2 TEXTBOOK FOR EACH SUBJECT

As discussed in Section 3.2, we design a textbook-guided example annotation pipeline to encompass both the *breadth of knowledge* and the *depth of reasoning*. The textbooks used for each subject are detailed in the following tables. They are selected by expert annotators and are recognized as authoritative references in their respective fields.

Subject	Textbook
Astronomy	1. <i>Foundations of Astrophysics</i> (Ryden & Peterson, 2020) 2. <i>Stellar Structure And Evolution</i> (Pols, 2011)
Biology	1. <i>Biology, 2nd Edition</i> (Clark et al., 2018a) 2. <i>Introduction to Agricultural Engineering Technology: A Problem Solving Approach, 4th Edition</i> (Field & Long, 2018) 3. <i>Introduction to Environmental Engineering, 5th Edition</i> (Davis & Cornwell, 2012) 4. <i>The Economy of Nature, 7th Edition</i> (Ricklefs, 2013) 5. <i>The Molecular Biology of the Cell, 6th Edition</i> (Alberts et al., 2014)
Chemistry	1. <i>Atkins' Physical Chemistry, 12th Edition</i> (Atkins et al., 2023) 2. <i>Chemistry, 2nd Edition</i> (Flowers et al., 2019) 3. <i>Chemistry: The Central Science, 15th Edition</i> (Brown et al., 2023) 4. <i>Organic Chemistry As A Second Language</i> (Klein, 2024) 5. <i>Organic Chemistry, 2nd Edition</i> (Clayden et al., 2012)
Electromagnetism	1. <i>Introduction to Electrodynamics, 4th Edition</i> (Griffiths, 2023) 2. <i>University Physics Volume 2 (Electromagnetism)</i> (Ling et al., 2016b)
Geography	1. <i>Fundamentals of Geophysics, 2nd Edition</i> (Lowrie & Fichtner, 2020) 2. <i>Human Geography, 12th Edition</i> (Fouberg & Murphy, 2020) 3. <i>Physical Geography: A Landscape Appreciation, 10th Edition</i> (Hess & McKnight, 2021)
Mechanics	1. <i>University Physics Volume 1</i> (Ling et al., 2016a)
Modern Physics	1. <i>University Physics Volume 3</i> (Ling et al., 2016c)
Neurobiology	1. <i>Neuroscience, 6th Edition</i> (Purves et al., 2018) 2. <i>Principles of Neural Science, 6th Edition</i> (Kandel et al., 2021) 3. <i>Principles of Neurobiology</i> (Luo, 2020)
Thermodynamics	1. <i>An Introduction to Thermal Physics</i> (Schroeder, 2020) 2. <i>University Physics Volume 2 (Thermodynamics)</i> (Ling et al., 2016b)

Table 6: List of textbooks and corresponding example numbers for the **Science** discipline.

Subject	Textbook
Biomedical Engineering	<ol style="list-style-type: none"> <li>1. <i>Biomaterials Science: An Introduction to Materials in Medicine, 4th Edition</i> (Wagner et al., 2020)</li> <li>2. <i>Biomaterials and Biopolymers</i> (Domb et al., 2023)</li> <li>3. <i>Fundamentals and Advances in Medical Biotechnology</i> (Anwar et al., 2022)</li> <li>4. <i>Introduction to Biomedical Engineering, 4th Edition</i> (Enderle &amp; Bronzino, 2017)</li> </ol>
Civil Engineering	<ol style="list-style-type: none"> <li>1. <i>Engineering Geology and Construction</i> (Bell, 2004)</li> <li>2. <i>Principles of Geotechnical Engineering, 9th Edition</i> (Das, 2017)</li> <li>3. <i>Structure for Architects: A Case Study in Steel, Wood, and Reinforced Concrete Design</i> (Bedi &amp; Dabby, 2019)</li> </ol>
Computer Science	<ol style="list-style-type: none"> <li>1. <i>Algorithms, 4th Edition</i> (Sedgewick &amp; Wayne, 2011)</li> <li>2. <i>Computer Organization and Design: The Hardware/Software Interface, 6th Edition</i> (Patterson &amp; Hennessy, 2022)</li> <li>3. <i>Computer Systems: A Programmer's Perspective, 3rd Edition</i> (Bryant &amp; O'Hallaron, 2011)</li> <li>4. <i>Deep Learning</i> (Goodfellow et al., 2016)</li> <li>5. <i>Digital Image Processing, 4th Edition</i> (Rafael &amp; Richard, 2018)</li> <li>6. <i>Introduction to Algorithms, 4th Edition</i> (Cormen et al., 2022)</li> <li>7. <i>Operating System Concepts, 10th Edition</i> (Silberschatz et al., 2018)</li> </ol>
Electrical Engineering	<ol style="list-style-type: none"> <li>1. <i>Electrical Engineering: Principles and Applications, 7th Edition</i> (Hambrey, 2018)</li> </ol>
Electronics and Communication	<ol style="list-style-type: none"> <li>1. <i>CMOS Analog Circuit Design, 3rd Edition</i> (Allen &amp; Holberg, 2011)</li> <li>2. <i>Introduction to Communication Systems</i> (Madhow, 2014)</li> <li>3. <i>The Art of Electronics, 3rd Edition</i> (Horowitz &amp; Hill, 2015)</li> </ol>
Materials Science	<ol style="list-style-type: none"> <li>1. <i>Composite Materials: Science and Engineering, 3rd Edition</i> (Chawla, 2012)</li> <li>2. <i>Convection in Porous Media, 5th Edition</i> (Nield &amp; Bejan, 2017)</li> <li>3. <i>Fiber-Reinforced Composites Materials, Manufacturing, and Design, 3rd Edition</i> (Mallick, 2007)</li> <li>4. <i>Materials Science and Engineering: An Introduction, 10th Edition</i> (Callister Jr &amp; Rethwisch, 2020)</li> </ol>
Mechanical Engineering	<ol style="list-style-type: none"> <li>1. <i>Industrial Automation: An Engineering Approach</i></li> <li>2. <i>Industrial Robotics Control: Mathematical Models, Software Architecture, and Electronics Design</i> (Frigeni, 2022)</li> <li>3. <i>Intelligent Manufacturing System and Intelligent Workshop</i> (Wang)</li> <li>4. <i>Machine Tool Practices, 11th Edition</i> (Kibbe et al., 2019)</li> <li>5. <i>Marks' Standard Handbook for Mechanical Engineers, 12th Edition</i> (Avalone et al., 2018)</li> <li>6. <i>Modern Control Engineering, 5th Edition</i> (Ogata, 2010)</li> </ol>

Table 7: List of textbooks and corresponding example numbers for the **Engineering** discipline.

Subject	Textbook
Basic Medicine	1. <i>Kuby Immunology, 8th Edition</i> (Owen et al., 2018) 2. <i>Robbins and Cotran Pathologic Basis of Disease, 10th Edition</i> (Kumar et al., 2020) 3. <i>Tissue Barriers in Disease, Injury and Regeneration</i> (Gorbunov, 2022)
Clinical Medicine	1. <i>Cecil Essentials of Medicine, 10th Edition</i> (Wing & Schiffman, 2021) 2. <i>Kumar and Clark's Clinical Medicine, 10th Edition</i> (Feather et al., 2020)
Dentistry	1. <i>Pharmacology and Therapeutics for Dentistry, 7th Edition</i> (Yagiela et al., 2010)
Pharmacy	1. <i>The Pharmacological Basis of Therapeutics, 13th Edition</i> (Brunton et al., 2017)
Preventive Medicine	1. <i>Public Health and Preventive Medicine, 15th Edition</i> (Maxcy et al., 2008)

Table 8: List of textbooks and corresponding example numbers for the **Healthcare** discipline.

Subject	Textbook
Art	1. <i>Art Through the Ages: A Global History Volume I, 16th Edition</i> (Kleiner, 2020) 2. <i>Introduction to Film Studies, 5th Edition</i> (Nelmes, 2012) 3. <i>The Filmmaker's Handbook: A Comprehensive Guide for the Digital Age, 5th Edition</i> (Ascher & Pincus, 2012)
Economics	1. <i>Intermediate Microeconomics: A Modern Approach, 8th Edition</i> (Varian, 2010) 2. <i>Land Resource Economics and Sustainable Development: Economic Policies and the Common Good</i> (Van Kooten, 2011) 3. <i>Macroeconomics, 9th Edition</i> (Blanchard, 2024) 4. <i>Principles of Economics, 3rd Edition</i> (Greenlaw et al., 2023) 5. <i>Principles of Microeconomics, 9th Edition</i> (Mankiw, 2020)
History	1. <i>Archaeology: Theories Methods and Practice, 7th Edition</i> (Renfrew & Bahn, 2016) 2. <i>World History Volume I: to 1500</i> (Kordas et al., 2022)
Law	1. <i>Arbitration Awards: A Practical Approach</i> (Turner, 2008) 2. <i>Contract Law</i> (Turner, 2013) 3. <i>The CISG: A new textbook for students and practitioners</i> (Huber & Mullis, 2009)
Literature	1. <i>An Introduction to Language, 11th Edition</i> (Fromkin et al., 2017) 2. <i>The Cambridge Introduction to the Novel</i> (MacKay, 2010)
Management	1. <i>Principles of Management</i> (Bright et al., 2019)

Table 9: List of textbooks and corresponding example numbers for the **Humanities and Social Science** discipline.

### A.3 ANNOTATION GUIDELINE AND INTERFACE

With the goal of ensure the high quality of data, **MMVU** adheres to the following four benchmark construction desiderata, we develop the following annotation interface based on Turkle (HLT-COE@JHU, 2024), an open-source clone of Amazon’s Mechanical Turk:

The screenshot shows the 'Video QA Annotation' interface. At the top, there's a navigation bar with 'Turkle Admin Stats Help' and a user status 'Logged in as [redacted] Change Password - Logout'. Below that is a project header 'Project: Annotation / Batch:'. On the right, there are buttons for 'Auto-accept next Task' (unchecked), 'Return Task' (red), 'Skip Task' (red), and 'Expires in 23:59'. The main area is titled 'Video QA Annotation' and contains fields: 'Enter a YouTube Video URL:' with the value 'https://www.youtube.com/watch?v=vp50Qggk4', 'Select the question type:' with a dropdown menu showing 'Multiple Choice', and a 'Submit Video' button. A small error message at the bottom says 'Error: The video does not have a Creative Commons license.' Below these fields is a 'Submit' button.

Figure 6: **Annotation Interface - Step 1: Video Collection.** In this step, annotators are required to input the YouTube video URL and select the desired question type. The backend system of the interface will automatically verify whether the provided YouTube video is under a Creative Commons license using the YouTube Data API v3. If the video does not meet this requirement, as shown in the figure, a warning message will be displayed, and the submission will be blocked. Once a valid example is submitted, the annotation interface will proceed to Step 2, which is illustrated in the following two figures.

The screenshot shows the 'Multiple-choice Question Annotation' interface. At the top, it has the same navigation and project header as Figure 6. The main area includes a video player showing a person pouring liquid into a beaker, with a timestamp '0:26 / 1:03'. Below the video are 'Start Time:' and 'End Time:' fields set to '00:07' and '00:55' respectively. A 'Question:' field contains the text: 'Assume that 2.24 liters of gas fully participates in the reaction shown in the video under the standard temperature and pressure condition, how many grams of precipitate are produced approximately?'. Under 'Textbook and Chapter:', there's a dropdown 'Chemistry, 2nd Edition (Paul Flowers, Klaus Theopold, Richard Langley)' and a search bar with '3'. A 'Knowledge (Wikipedia)' section lists 'Calcium hydroxide', 'Carbon dioxide', and 'Ideal gas law' with their corresponding Wikipedia URLs. A 'Reasoning Process:' section contains two numbered points: '1. In the video, a person exhales gas that is continuously introduced into a clear solution, gradually forming a white precipitate. This indicates that the substances involved in the reaction are CO<sub>2</sub> and limestone.' and '2. The chemical reaction equation is: Ca(OH)<sub>2</sub> + CO<sub>2</sub> → CaCO<sub>3</sub> + H<sub>2</sub>O'. At the bottom is a 'Submit' button.

Figure 7: Annotation Interface - Step 2: Multiple-choice Question Annotation.

Turkle Admin Stats Help

Logged in as [ ] Change Password - Logout

Project: Annotation / Batch:

Auto-accept next Task   Expires in 23:55

### Video QA Annotation

Question:  
Which cinematic shooting technique is shown in the video?

Answer:  
Dolly Zoom

Start Time: 00:00 End Time: 00:15

Textbook and Chapter:  
Introduction to Film Studies, 5th Edition (Jill Nelmes) 4

Knowledge (Wikipedia):

- Dolly zoom
- Cinematic techniques
- + Add Wikipedia Page

Reasoning Process:  
"Dolly zoom" is a unique cinematographic technique that creates a visual distortion effect by simultaneously adjusting the lens while dollying in or out. The video clearly shows the camera zooming in on the object. Therefore the most possible shooting technique demonstrated in the video is dolly zoom.

Figure 8: Annotation Interface - Step 2: Open-ended Question Annotation.

#### A.4 VALIDATION GUIDELINE AND INTERFACE

To ensure that the final dataset remains high-quality and meets expert-level standards without introducing unnecessary bias, each example in **MMVU** undergoes expert review by one of the authors or top-performing annotators to verify the accuracy of its annotations, following the annotation guideline detailed in Appendix A.3. The examples of validation interface are presented as follows:

The screenshot shows the MMVU validation interface. At the top, there's a navigation bar with 'Turkle', 'Admin', 'Stats', 'Help', and a user status 'Logged in as [redacted] | Change Password - Logout'. Below the navigation is a header 'Project: Validation / Batch: | VIDEO QA Validation'. On the left, a video player shows a person in a lab setting. The video progress is at 0:08 / 0:46. Below the video is a checked checkbox labeled 'Checked'. To the right of the video is a question: 'Assume that 2.24 liters of gas fully participates in the reaction shown in the video under the standard temperature and pressure condition.' Below the question are five numerical options: 10.0, 5.0, 12.0, 15.0, and 20.0, each with a 'Checked' checkbox next to it. Below the options is a dropdown menu labeled 'Select the Correct Answer:  Checked' with the value 'A'. Underneath this is a section titled 'Knowledge (Wikipedia):  Checked' with three links: 'Calcium hydroxide' ([https://en.wikipedia.org/wiki/Calcium\\_hydroxide](https://en.wikipedia.org/wiki/Calcium_hydroxide)), 'Carbon dioxide' ([https://en.wikipedia.org/wiki/Carbon\\_dioxide](https://en.wikipedia.org/wiki/Carbon_dioxide)), and 'Ideal gas law' ([https://en.wikipedia.org/wiki/Ideal\\_gas\\_law](https://en.wikipedia.org/wiki/Ideal_gas_law)). Below this is a section titled 'Reasoning Process:  Checked' with two numbered points: '1. In the video, a person exhales gas that is continuously introduced into a clear solution, gradually forming a white precipitate. This indicates that the substances involved in the reaction are CO<sub>2</sub> and limewater.' and '2. The chemical reaction equation is: Ca(OH)<sub>2</sub> + CO<sub>2</sub> → CaCO<sub>3</sub> + H<sub>2</sub>O'. At the bottom of the interface are three buttons: 'Submit' (blue), 'Require the annotator to revise' (yellow), and 'Discard example due to low quality' (red).

**Figure 9: Validation Interface.** Human validators are required to thoroughly review each annotation feature to ensure alignment with benchmark construction criteria and annotation guidelines. If revisions are not feasible, detailed feedback must be provided to the original annotator, who will then revise and resubmit the annotation for a second review. Additionally, validators may discard examples deemed to be of low quality and unlikely to meet the desired criteria through revision.

#### A.5 DATA ANNOTATION AND VALIDATION PAYMENT

The annotation and validation process for **MMVU** spans three months. As outlined in Section 3.2, annotating examples for **MMVU** can be particularly time-intensive, especially when there is limited availability of videos with Creative Commons licenses in the required subjects. To accommodate this and ensure a high-quality dataset, we compensate annotators based on the time they spend rather than the number of examples completed, preventing them from rushing through tasks. Annotators are required to record their screens throughout the annotation process, which enables us to verify time reporting accuracy and maintain productivity standards. This also helps us identify any distractions and precisely track the total time spent on each task. We offer a *base rate* of 6 USD per hour for both annotation and validation work, with an additional 2 USD per completed annotation and 0.40 USD per validated example. On average, annotating a single question for **MMVU** takes 20 minutes and 17 seconds, while validation requires 4 minutes and 12 seconds. This compensation structure ensures that annotators earn wages that are competitive with the average payment for teaching assistants at their respective universities. To reduce pressure and maintain a comfortable pace, we recommended that annotators limit their work to a maximum of 10 QA example annotations or 50 QA example validations per day.

## B EXPERIMENT SETUP

### B.1 CONFIGURATION OF EVALUATED MODELS

Table 10 detail the configuration of each evaluated models. We use the default settings from the official implementation of each model to process vision input. Across all experiments, the temperature is set to 1.0, with a maximum output length of 1024 tokens. However, for Gemini-2-Flash-Thinking, the maximum output length is set as 8192 tokens to accommodate its long CoT reasoning mechanism. All inferences are reproducible on a workstation equipped with two NVIDIA A100-80G GPUs.

Organization	Model	Release	Version	Support Video?	Input Frames	# Inference Pipeline
<i>Proprietary Models</i>						
OpenAI	o1*	2024-12	o1-2024-12-17	✗	32	API
	GPT-4o	2024-8	gpt-4o-2024-08-06	✗	32	
	GPT-4o-mini	2024-7	gpt-4o-mini-2024-07-18	✗	32	
Google	Gemini 2.0 Flash Thinking	2024-12	gemini-2.0-flash-thinking-exp-1219	✗	32	API
	Gemini 2.0 Flash	2024-12	gemini-2.0-flash-exp	✗	32	
	Gemini 1.5 Pro	2024-9	gemini-1.5-pro	✓	32	
	Gemini 1.5 Flash	2024-9	gemini-1.5-flash	✓	32	
Anthropic	Claude-3.5-Sonnet	2024-10	claude-3-5-sonnet-20241022	✗	32	API
xAI	Grok-2-Vision	2024-12	grok-2-vision-1212	✗	32	API
Zhipu AI	GLM-4V-Plus	2025-1	glm-4v-plus-0111	✓	4	API
<i>Open-source Multimodal Foundation Models</i>						
Mistral AI	Pixtral-12B	2024-9	Pixtral-12B-2409	✗	8	vLLM
Microsoft	Phi-3.5-Vision	2024-7	Phi-3.5-vision-instruct	✗	16	vLLM
Shanghai AI Lab	InternVL2.5-38B	2024-11	InternVL2.5-38B	✗	4	vLLM
	InternVL2.5-8B	2024-11	InternVL2.5-8B	✗	4	
	InternVL2-8B	2024-6	InternVL2-8B	✗	4	
Alibaba	Qwen2-VL-2B	2024-8	Qwen2-VL-2B-Instruct	✓	1fps	vLLM
	Qwen2-VL-7B	2024-8	Qwen2-VL-7B-Instruct	✓	1fps	
	Qwen2-VL-72B	2024-9	Qwen2-VL-72B-Instruct	✓	1fps	
Meta	Llama-3.2-11B-Vision	2024-9	Llama-3.2-11B-Vision-Instruct	✗	8	vLLM
	Llama-3.2-90B-Vision	2024-9	Llama-3.2-90B-Vision-Instruct	✗	8	
DAMO	VideoLLaMA2-7B	2024-6	VideoLLaMA2-7B	✓	1fps	HF
	VideoLLaMA2.1-7B	2024-10	VideoLLaMA2.1-7B-16F	✓	1fps	
DeepSeek	DeepSeek-VL2	2024-12	deepseek-vl2	✗	2	vLLM
	DeepSeek-VL2-Small	2024-12	deepseek-vl2-small	✗	2	
	DeepSeek-VL2-Tiny	2024-12	deepseek-vl2-tiny	✗	2	
Rhymes	Aria	2024-11	Aria-Chat	✗	8	vLLM
Llava Hugging Face	Llava-OneVision-7B	2024-9	llava-onevision-qwen2-7b-ov-chat-hf	✓	1fps	vLLM
	Llava-NeXT-Video-34B	2024-6	Llava-NeXT-Video-34B-hf	✗	8	
	Llava-NeXT-Video-7B	2024-6	Llava-NeXT-Video-7B-hf	✓	16	
HuggingFaceM4	Idefics3-8B	2024-8	Idefics3-8B-Llama3	✗	4	vLLM
OpenGVLab	InternVideo2-8B	2024-8	InternVideo2-Chat-8B	✓	1fps	HF
H2O	H2OVL Mississippi-2B	2024-10	h2ovl-mississippi-2b	✗	4	vLLM

Table 10: Details of the multimodal foundation models evaluated in **MMVU**. The “Source” column includes URLs for proprietary models and Hugging Face model names for open-source models. The “# Input Frames” column, for those models only support multi-image input, represents the default number of input frames, chosen from 2, 4, 8, 16, 32, based on the maximum value that does not exceed the model’s context window. “HF” means “Hugging Face”.

## B.2 CHAIN-OF-THOUGHT AND DIRECT ANSWER PROMPTS

The following figures illustrates the CoT reasoning and Direct Answer prompts applied in this study for answering multiple-choice and open-ended questions, respectively.

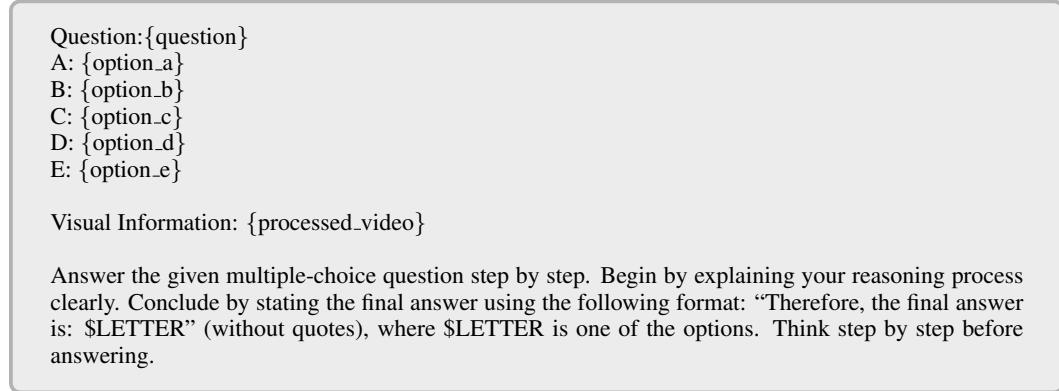


Figure 10: CoT reasoning prompt, adopted from MMMU-Pro (Yue et al., 2024b), for answering multiple-choice question.

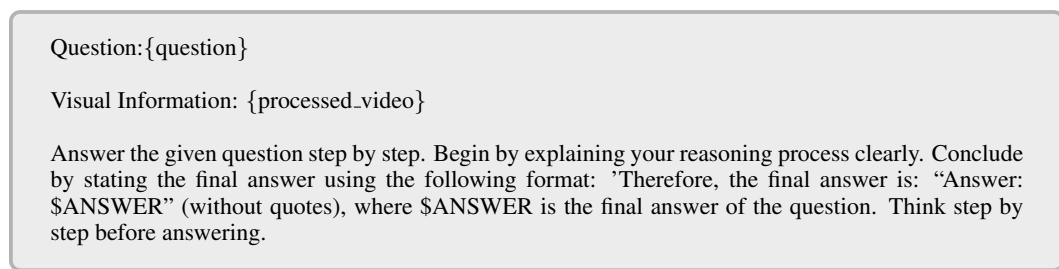


Figure 11: CoT reasoning prompt for answering open-ended question.

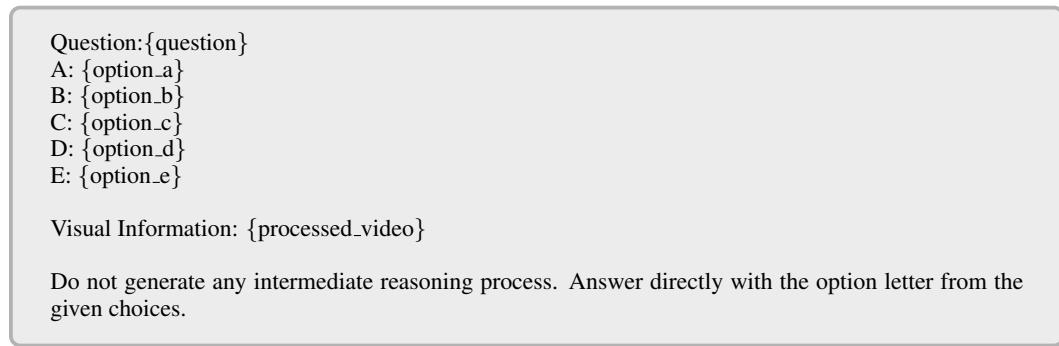


Figure 12: Direct Answer prompt, adopted from MMMU-Pro (Yue et al., 2024b), for answering multiple-choice question.

Question:{question}

Visual Information: {processed\_video}

Do not generate any intermediate reasoning process. Directly output the final answer.

Figure 13: Direct Answer prompt for answering open-ended question.

### B.3 PROMPTS FOR ACCURACY EVALUATION

[Instruction]

Evaluate whether the model’s final answer is correct by comparing it to the ground-truth answer provided for the given question.

You should first extract the final answer from the model’s response, and then compare the extracted answer with the ground-truth answer to determine its accuracy. Output your response in the following structured format:

```
{  
    'extracted_answer': // str value "A" "B" "C" "D" "E", should be a single character  
    'correct': // boolean value, True if the answer is correct, False otherwise  
}
```

[User]

Question:{question}

A: {option\_a}  
B: {option\_b}  
C: {option\_c}  
D: {option\_d}  
E: {option\_e}

Ground Truth Answer: {ground\_truth}

Model Response to the Question: {model\_response}

Figure 14: Evaluation prompt used for assessing the accuracy of multi-choice QA.

[Instruction]

Evaluate whether the model’s final answer is correct by comparing it to the ground-truth answer provided for the given question. You should first extract the final answer from the model’s response, and then compare the extracted answer with the ground-truth answer to determine its accuracy. The final answer generated by the model does not need to match the ground-truth answer word-for-word. However, it should only be considered correct if it demonstrates the exact same technique or concept explicitly and unambiguously equivalent to the ground-truth answer. Output your response in the following structured format:

```
{  
    'extracted_answer': // str value, the short final answer extracted from the model's response, do not  
    // hallucinate one that is not present in the response  
    'correct': // boolean value, True if the answer is correct, False otherwise  
}
```

[User]

Question:{question}

Ground Truth Answer: {ground\_truth}

Model Response to the Question: {model\_response}

Figure 15: Evaluation prompt used for assessing the accuracy of open-ended QA.

## C EXPERIMENT

### C.1 COMPARISON BETWEEN COT REASONING AND DIRECT ANSWERING

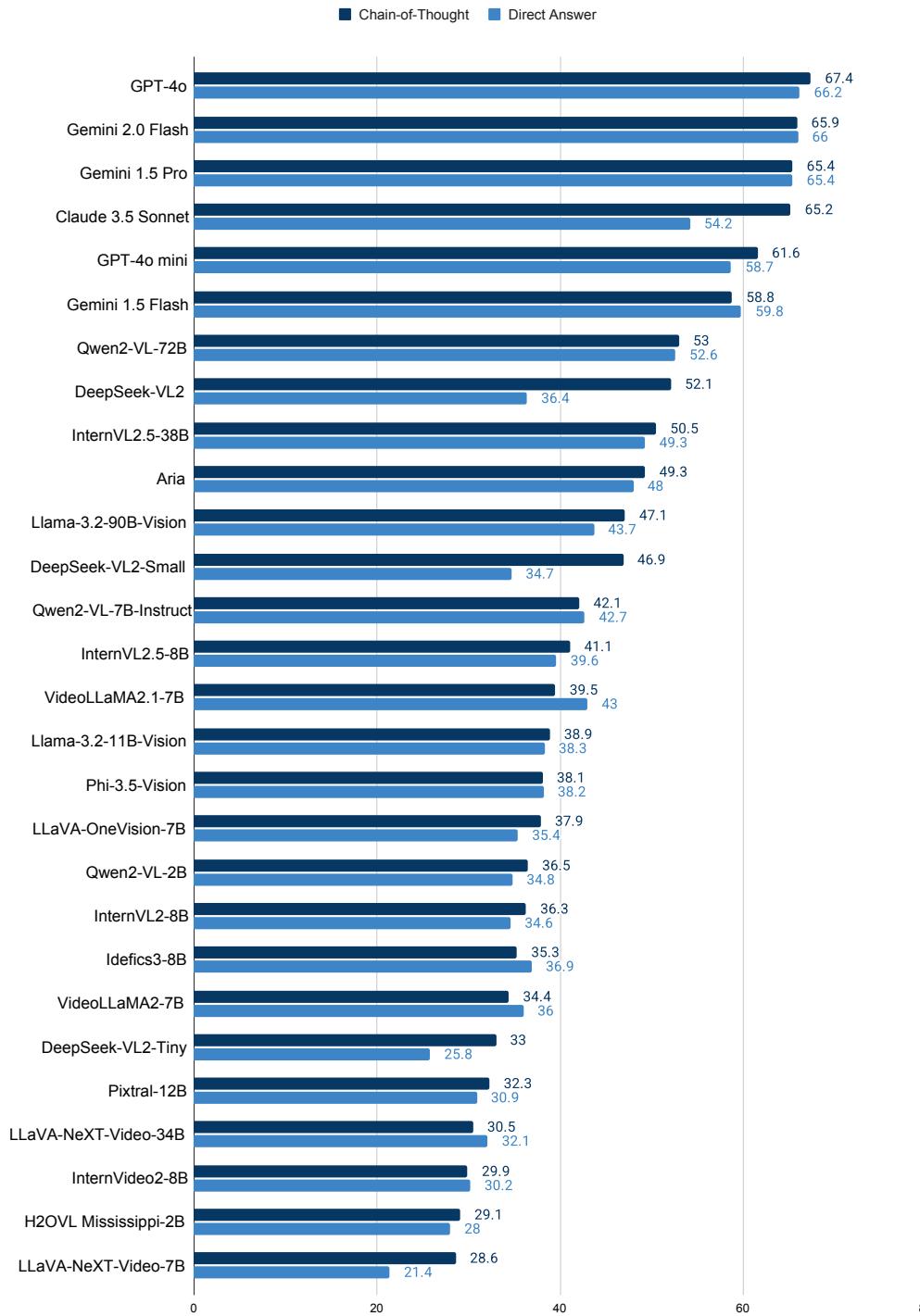


Figure 16: Comparison of model performance between CoT reasoning and direct answering on the validation set.

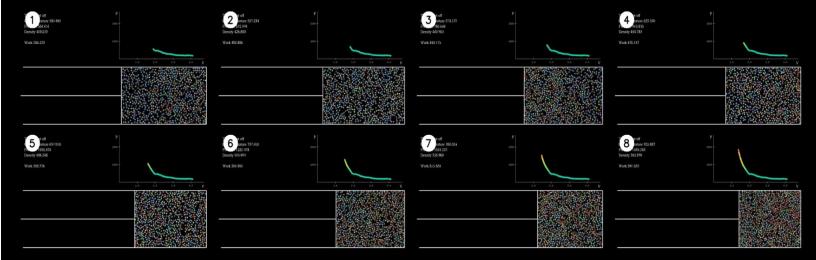
## C.2 ERROR CASE ANALYSIS: VISUAL PERCEPTION ERROR

**Thermodynamics**

**Question:** Which process is shown in the animation?

**Options:**

(A) Isothermal compression      (B) Adiabatic compression      (C) Isobaric compression  
 (D) Adiabatic expansion      (E) Isothermal expansion



**Answer:** B

**Relevant Knowledge (Wikipedia page title):** Adiabatic process, Isothermal process, Isobaric process, Ideal gas law

**Reasoning Rational:**  
 From the video, it can be seen that the gas is thermally isolated and returns to its initial state through compression. This physical process is known as adiabatic compression.

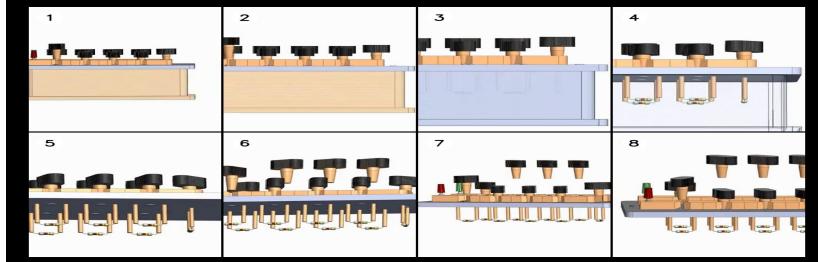
**Llama-3.2-90B-Vision:**  
 To determine the process shown in the animation, we need to analyze the provided graph and the printed values.  
 The graph plots pressure (P) against volume (V), and the curve is descending, indicating that the volume is increasing while the pressure decreases. This is characteristic of expansion.  
 Additionally, the absence of a thermostat indicates that no temperature regulation is involved. This suggests the process is not isothermal (constant temperature), but rather adiabatic (no heat exchange with the surroundings).  
 Since the volume is increasing, it is likely that the gas is expanding. (Human Evaluator's Comment: This observation contradicts with the video content, where the gas is being compressed.) Adiabatic expansion occurs when a gas expands without gaining or losing heat, and its temperature decreases.  
 Therefore, the final answer is: D. Adiabatic expansion.

**Error Category:** Visual perception error  
**Error Reason:** The model wrongly identifies the process of gas compression as the process of gas expansion, which is not consistent with the human observation and leads to the wrong final answer.

Figure 17: An error case of Thermodynamics.

## Electromagnetism

**Question:** Which physical variable is changed in the animation?



**Answer:** Resistance

**Relevant Knowledge (Wikipedia page title):** Electrical resistance and conductance, Ohm's law,

Electromagnetism

**Reasoning Rational:**

The video shows that by pulling out the knob, the resistance in the circuit can be changed.

**DeepSeek-VL2:**

In the animation, a PSP (Pressure Sensor Package) is shown being compressed by water, and its output increases as the depth increases due to the increased pressure and density of the water. ([Human Evaluator's Comment: This observation contradicts with the video content, where there is no water.](#)) This indicates that the deformation variable is the one that changes in the animation.

Therefore, the final answer is: Deformation

**Error Category:** Visual Perception Error

**Error Reason:** The model incorrectly identifies the device shell, which gradually turns blue and transparent in the video for inside demonstration, as water. The hallucination misleads the reasoning process and results in a totally wrong conclusion.

Figure 18: An error case of Electromagnetism.

## Art

**Question:** Which cinematic shooting technique is shown in the video?



**Answer:** Dolly Zoom

**Relevant Knowledge (Wikipedia page title):** Dolly zoom, Cinematic techniques

**Reasoning Rational:**

"Dolly zoom" is a unique cinematographic technique that creates a visual distortion effect by simultaneously adjusting the lens while dollying in or out. The video clearly shows the camera zooming in on the object. Therefore, the most likely shooting technique demonstrated in the video is the dolly zoom.

**Qwen2-VL-72B:**

The video shows the panning technique, where the camera moves horizontally, from left to right gradually. ([Human Evaluator's Comment: This observation contradicts with the video content, where the camera zooming in on the object, instead of moving horizontally.](#)) to capture a wide view of the scene.

Therefore, the answer is: panning.

**Error Category:** Visual Perception Error

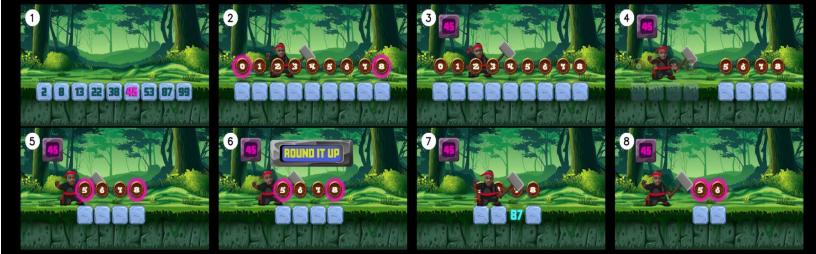
**Error Reason:** The model fails to interpret the dolling-in motion in the video or associate it with a cinematic technique. Instead, it thinks the motion of camera is panning from left to right. Therefore, the model fails to effectively interpret temporal visual data and it infers a wrong answer.

Figure 19: An error case of Art.

### C.3 ERROR CASE ANALYSIS: MISUSE OR LACK DOMAIN KNOWLEDGE IN VISUAL PERCEPTION

**Computer Science**

**Question:** What is the algorithm shown in the video?



The screenshots show a character interacting with a sequence of numbered blocks (1, 2, 3, 4, 5, 6, 7, 8) in a forest setting. The sequence starts at the top left and moves right, then down to the bottom left. In each frame, the character is either pointing at or moving one of the blocks. Frame 1 shows the full sequence. Frame 2 shows the character pointing at block 4. Frame 3 shows the character pointing at block 7. Frame 4 shows the character pointing at block 8. Frame 5 shows the character pointing at block 5. Frame 6 shows the character pointing at block 6 with a 'ROUND IT UP' button visible. Frame 7 shows the character pointing at block 7. Frame 8 shows the character pointing at block 8.

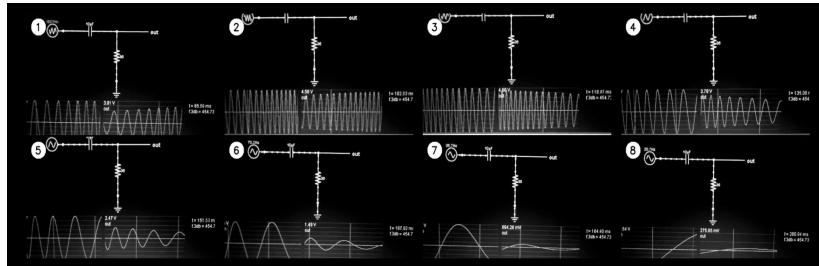
**Answer:** Binary Search  
**Relevant Knowledge (Wikipedia page title):** [Binary search](#), [Search algorithm](#), [Computer Science](#)  
**Reasoning Rational:**  
 The character in the video always locates data from the middle and then discards a part of the data, ultimately finding the target value. From the entire process, we can infer that this sequence is an ordered sequence, which aligns with the process of Binary Search.

**Llama-3.2-90B-Vision:**  
 To determine what algorithm is shown in the video, let's analyze the steps and their characteristics:  
 1. The video starts by showing a set of numbers.  
 2. It then rearranges the numbers in a specific order.  
 The observed steps suggest that the algorithm is [sorting the numbers in ascending order](#). The algorithm sorts an array by repeatedly finding the minimum element from the unsorted part and putting it at the beginning. It involves selecting the smallest element and moving it to the correct position, which matches the steps shown in the video.  
 Based on the analysis, the algorithm shown in the video involves [selecting the minimum element and moving it to the correct position repeatedly](#), which is characteristic of the selection sort algorithm. ([Human Evaluator's Comment: Misuse of domain knowledge in visual perception. The indices of the array are mistakenly identified as the values.](#))  
 Therefore, the final answer is: Selection Sort.

**Error Category: Misuse or Lack Domain Knowledge in Visual Perception**  
**Error Reason:** The model wrongly recognized the numerical sequence in the video as values (when in fact they are the indices of the sequence), leading to the incorrect conclusion that the algorithm shown in the video is a sorting algorithm.

Figure 20: An error case of Computer Science.

**Question:** What kind of filter is the circuit implemented?



**Answer:** High-pass filter

**Relevant Knowledge (Wikipedia page title):** High-pass filter, RC circuit, Electronic filter, Signal processing

**Reasoning Rational:**

From the video, it can be seen that the circuit processes different frequency components of the signal, filtering out low frequencies to produce a new signal. Thus, it is a High-pass filter.

#### Qwen2-VL-72B:

To determine the type of filter implemented in the circuit, we need to analyze the components and their configuration.

1. Inductor and Capacitor Configuration:

- The circuit includes an inductor and a capacitor.
- These components are arranged in a series configuration.

2. Circuit Analysis:

- The circuit diagram shows a series configuration of an inductor and a capacitor. ([Human Evaluator's Comment: Misuse of domain knowledge in visual perception. The resistor is mistakenly identified as an inductor.](#))
- This configuration is characteristic of a band-pass filter, where the inductor passes high frequencies and the capacitor passes low frequencies, but together they allow a specific range of frequencies to pass.

3. Confirmation:

- The given figure and the analysis of the components and their arrangement strongly suggest that this is a band-pass filter.

Therefore, the final answer is: Band-pass filter.

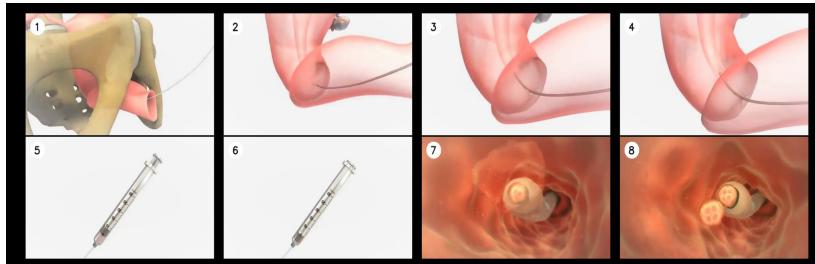
**Error Category: Misuse or Lack Domain Knowledge in Visual Perception**

**Error Reason:** The model lacks domain knowledge of common circuit component diagrams, leading to the incorrect identification of the resistor in the circuit as an inductor. The wrong identification of circuit components results in the false reasoning of the filter type.

Figure 21: An error case of Electrical Engineering.

## Pharmacy

**Question:** What process is shown in the video?



**Answer:** Embryo Transfer

**Relevant Knowledge (Wikipedia page title):** Embryo transfer, In vitro fertilisation, Childbirth, Female reproductive system

**Reasoning Rationale:**

The video demonstrates the process of embryo transfer in which an embryo is implanted into a woman's uterus following in vitro fertilization. It features specific visuals of the fertilized egg, syringe, and the injection procedure.

**DeepSeek-VL2:**

The video shows a syringe injecting a cell into a tube, leading to the creation of baby DNA and cells multiplying, eventually forming a dark blue circular mass with an umbilical cord attached to it. This process is illustrated with illustrations of babies inside, suggesting fetal development. (*Human Evaluator's Comment: Misuse of domain knowledge in visual perception. The uterus of a woman where the process happens is mistakenly identified as the inside of a fetus.*) Therefore, the final answer is: Fetal Development Process.

**Error Category: Misuse or Lack Domain Knowledge in Visual Perception**

**Error Reason:** The model misuses its domain knowledge of fetal development, wrongly identifying the uterus shown in the video as a fetus. Then the model draws the incorrect conclusion that the video demonstrates the fetal development process.

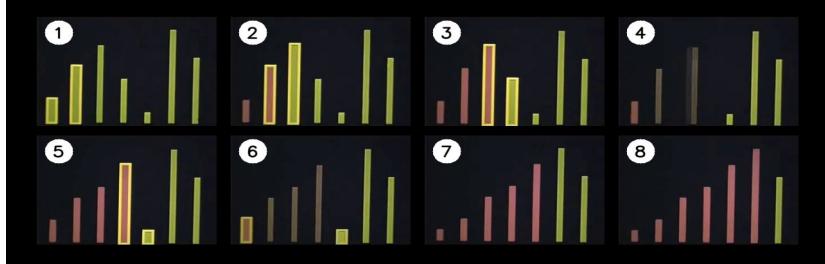
Figure 22: An error case of Pharmacy.

#### C.4 ERROR CASE ANALYSIS: MISUSE OR LACK DOMAIN KNOWLEDGE IN REASONING

**Computer Science**

**Question:** The algorithm demonstrated in the video is which common algorithm?

**Options:**  
 (A) Bubble Sort (B) Merge Sort (C) Insertion Sort (D) Selection Sort (E) Heap Sort



**Answer:** C

**Relevant Knowledge (Wikipedia page title):** [Insertion sort](#), [Bubble sort](#), [Merge sort](#), [Selection sort](#), [Heapsort](#)

**Reasoning Rational:**  
 The video ultimately sorts these blocks in ascending order, thus it is clearly a sorting algorithm. Observing the video process, there is no need to scan the sorting process from beginning to end, and there is an obvious insertion action, which allows us to eliminate Selection Sort, indicating it is Insertion Sort.

**DeepSeek-VL2:**  
 The given video shows computer bars that increase and decrease according to their size. The algorithm being demonstrated is a sorting algorithm, as it reorganizes the bars in order of their size. Among the given options, options (B) and (C) are not suitable, as [Merge Sort](#) and [Insertion Sort](#) are not typically demonstrated through visual bar representations. Option (F) is also not suitable, as [Heap Sort](#) is not typically demonstrated through visual bar representations. ([Human Evaluator's Comment: Biased misuse of pre-trained domain knowledge in reasoning. The sorting algorithms can be demonstrated through visual bar representations and the form of presentation should not be the main focus.](#)) Although option (D) Selection Sort is a possibility, the given algorithm does not repeatedly find the minimum element of the unsorted part and put it at the front, which characterizes Selection Sort. Therefore, the given algorithm is most likely a Bubble Sort, as this algorithm uses visual bar representations.

Therefore, the final answer is: A. Bubble Sort

**Error Category: Misuse or Lack Domain Knowledge in Reasoning**

**Error Reason:** The model considers Insertion Sort not to be the correct answer because it believes that Insertion Sort is not typically demonstrated through visual bar representations. This is a typical case of reasoning errors caused by the misuse of pretrained domain knowledge.

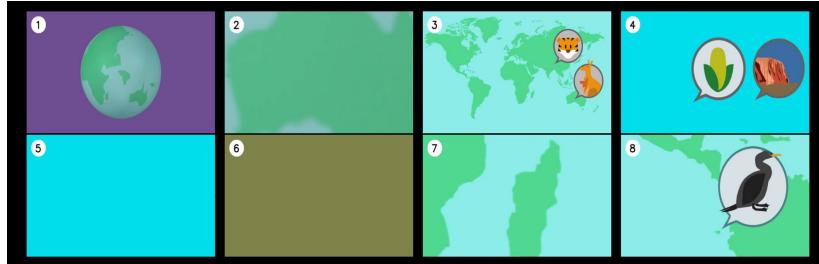
Figure 23: An error case of Computer Science.

## Biology

**Question:** What biological process is shown in the video?

**Options:**

- (A) Biodegradation (B) Community succession (C) Genetic drift (D) Migration of organism (E) Species formation



**Answer:** E

**Relevant Knowledge (Wikipedia page title):** Speciation, Reproductive isolation, Allopatric speciation, Charles Darwin, On the Origin of Species

**Reasoning Rational:**

In the video, the species undergo geographical isolation, which prevents gene flow with other populations. Over time, this leads to changes in allele frequencies within the gene pool. When significant differences arise between two populations' gene pools causing reproductive isolation, a new species is formed.

**Owen2-VL-72B:**

The video shows a girl floating over a virtual map and selecting various animals. This image is used to demonstrate the process of migration of organisms. The girl's selection of animals and their placement on the map suggests the migration of organisms to different locations. (*Human Evaluator's Comment: Misuse of domain knowledge in reasoning. The movement of organisms is due to the continental drift, not a spontaneous migration.*) Additionally, the video features a tree, a fish, a koala, a giraffe, and a lion on the edges of the map, which further supports the notion of migration. Based on this analysis, the biological process shown in the video is the migration of organisms.

Therefore, the final answer is: {D. Migration of organisms}

**Error Category: Misuse or Lack Domain Knowledge in Reasoning**

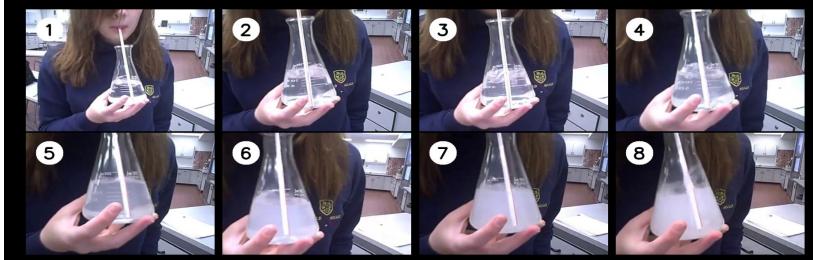
**Error Reason:** The model misuses pretrained knowledge and incorrectly infers that this is a case of spontaneous organism migration. In fact, the organism movement and isolation is caused by continental drift.

Figure 24: An error case of Biology.

## Chemistry

**Question:** Assume that 2.24 liters of gas fully participates in the reaction shown in the video under the standard temperature and pressure condition, how many grams of precipitate are produced approximately?

**Options:** (A) 10.0 (B) 5.0 (C) 12.0 (D) 15.0 (E) 20.0



**Answer:** A

**Relevant Knowledge (Wikipedia page title): Calcium hydroxide, Carbon dioxide, Ideal gas law**

**Reasoning Rational:**

1. In the video, a person exhales gas that is continuously introduced into a clear solution, gradually forming a white precipitate. This indicates that the substances involved in the reaction are CO<sub>2</sub> and limewater.
2. The chemical reaction equation is: Ca(OH)<sub>2</sub> + CO<sub>2</sub> → CaCO<sub>3</sub> + H<sub>2</sub>O
3. At the STP, 2.24 liters of CO<sub>2</sub> corresponds to 0.1 Moles.
4. From balanced equation, 0.1 moles of CO<sub>2</sub> produce 0.1 moles of CaCO<sub>3</sub>.
5. Given Ca = 40 g/mol, C = 12 g/mol, O = 16 g/mol, the molar mass of CaCO<sub>3</sub> = 40 + 12 + 16 \* 3 = 100 g/mol.  
Therefore, the mass of CaCO<sub>3</sub> = 0.1 \* 100 = 10g.

### DeepSeek-VL2:

To determine the amount of precipitate produced, we need to identify the chemical reaction involved and use the stoichiometry of the reaction. However, the video does not provide enough information to directly identify the reaction (**Human Evaluator's Comment: Lack of domain knowledge in reasoning. The model fails to recall and apply relevant chemical equations.**). Assuming the reaction is a typical precipitation reaction involving a gas and a solution, we can use the ideal gas law (PV=nRT) to estimate the volume of gas produced under standard temperature and pressure (STP).

Given:

- Volume of gas produced = 2.24 liters
- STP conditions: P = 1 atm, T = 273 K, R = 0.0821 L·atm/(mol·K)

First, we calculate the number of moles of gas produced using the ideal gas law:  
$$n = PV/RT = (1 \text{ atm} * 2.24 \text{ L}) / (0.0821 \text{ (L·atm)/(mol·K)} * 273 \text{ K}) \approx 0.1 \text{ mol}$$

Assuming the reaction produces a precipitate with a known stoichiometry, we can calculate the mass of the precipitate. However, without knowing the specific reaction, we cannot accurately determine the mass of the precipitate.

Given the options, we need to make an educated guess based on typical reactions. If we assume the reaction produces a precipitate with a molar mass of around 50 g/mol (a common value for many precipitates), we can estimate the mass of the precipitate as follows:

$$\text{Mass of precipitate} = n * \text{Molar mass} = 0.1 \text{ mol} * 50 \text{ g/mol} \approx 5 \text{ g}$$

Therefore, the correct answer is (B) 5.0 g.

### Error Category: Misuse or Lack Domain Knowledge in Reasoning

**Error Reason:** The model fails to recall and apply the correct relevant chemical equations. Instead, it merely makes a vague statement and then guesses a possible ratio between the products and reactants, which comes to a wrong conclusion at last.

Figure 25: An error case of Chemistry.

## C.5 ERROR CASE ANALYSIS: HEAVY RELIANCE ON TEXTUAL INFORMATION

**Clinical Medicine**

**Question:** When the phenomenon in the video occurs, which of the following options is the most important to adjust?

**Options:**

- (A) Renin-angiotensin system
- (B) Hypothalamic-pituitary-thyroid axis
- (C) Adrenaline and norepinephrine
- (D) Melatonin secretion
- (E) Sodium-potassium pump activity

**Answer:** A

**Relevant Knowledge (Wikipedia page title):** Renin–angiotensin system, Homeostasis, Assessment of kidney function, Osmoregulation

**Reasoning Rational:**  
The issue presented in the video is related to the imbalance of water within the body, causing cells to swell and resulting in symptoms such as headaches. Among the options provided, the regulation of water balance is primarily managed through the renin-angiotensin system.

**Llama-3.2-90B-Vision:**  
The phenomenon in the video likely refers to the effects of overhydration on the body. To adjust this, **the most important physiological response would be to regulate the body's sodium level, which is primarily managed by the sodium-potassium pump activity.** (**Human Evaluator's Comment: The model ignores the video information about water balance.**) This pump plays a crucial role in maintaining the balance of electrolytes and fluids within the body, which is crucial for the proper functioning of cells.

Therefore, the final answer is: {E. Sodium-potassium pump activity}.

**Error Category: Heavy Reliance on Textual Information**

**Error Reason:** The model ignores the video information about water excretion in the human body. It relies solely on the text and chooses the incorrect option E simply because cellular ion balance is crucial to the body's internal environment, thereby missing the actual disorder of fluid overload mentioned in the video.

Figure 26: An error case of Clinical Medicine.

## Management

**Question:** For the company shown in the video, what is one of the methods to handle associated risks?  
**Options:**

- (A) diversification (B) vertical integration (C) horizontal integration (D) outsourcing (E) speculation



**Answer:** B

**Relevant Knowledge (Wikipedia page title):** Vertical integration, Horizontal integration, Diversification, Risk management, Supply chain management

### Reasoning Rational:

According to the content of the video, the company being discussed is a coffee company (the bubbles above the truck represent coffee beans) along with its upstream supply chain (costs of raw coffee beans, labor costs, etc.). A common practice to mitigate the price risk of upstream products is vertical integration, which involves establishing the coffee brand's own coffee plantation.

### Qwen2-VL-72B:

To handle the risk associated with a company, diversification is one of the methods that can be used. Diversification involves spreading investments across various sectors or industries to reduce risk. Vertical integration and horizontal integration are strategies to gain control over the supply chain or expand market dominance, respectively, but they do not directly address risk management ([Human Evaluator's Comment: The model ignores the video information about the features of the company](#)). Inventory management and outsourcing are operational strategies that can influence risk, but they are not specifically methods for handling risk. Therefore, the correct answer is Option A. Therefore, the final answer is: A. diversification.

### Error Category: Heavy Reliance on Textual Information

**Error Reason:** The model ignores the features of the company demonstrated in the video. Instead, the model jumps to the conclusion that (A) diversification is the most appropriate method of risk management, solely based on the question and the given options.

Figure 27: An error case of Management.

## C.6 ERROR CASE ANALYSIS: LOGICAL REASONING ERROR

**Mechanical Engineering**

**Question:** What kind of mechanical transformation is not represented in the video?

**Options:**

- (A) It may convert gravitational energy into elastic energy.
- (B) It may convert one force to another force.
- (C) It may convert one speed to another speed.
- (D) It may convert linear motion into rotational motion.
- (E) It may convert gravitational energy to kinetic energy.

**Answer:** D

**Relevant Knowledge (Wikipedia page title):** Mechanical energy, Kinematics, Mechanical Engineering

**Reasoning Rational:**

In the video, the rotation of the disk carries the vertical movement of the animal block, and the speed of the block changes during the start and stop. The gravity of the mass and the elasticity of the spring transform each other. So the choice is D.

**Llama-3.2-90B-Vision:**

To determine which type of mechanical transformation is not represented in the video, we need to consider the actions depicted in the video and the various forms of energy and motion transformation mechanisms typically found in mechanical systems.

1. Linear Motion into Rotational Motion: The video shows a mechanical linkage where a linear motion is converted into rotational motion.
2. Rotational Motion into Linear Motion: The video also shows the reverse process, where rotational motion is converted into linear motion, which is visible in the reciprocal motion of the device.
3. Kinetic Energy Transformation: The video shows a transformation of kinetic energy from one part of the mechanism to another as parts move and interact.
4. Elastic Energy Transformation: The presence of a spring in the mechanism suggests the conversion of kinetic energy into elastic potential energy and vice versa.

Now, let's evaluate the options to see which transformation is not explicitly shown in the video:

Option A: It may convert gravitational energy into elastic energy. - **This transformation is not directly shown in the video. (Human Evaluator's Comment: The reasoning here contradicts with the other reasoning.)** The video focuses on the interactions between mechanical parts and their resulting motions rather than showing the specific conversion of gravitational energy into elastic energy.

Option B: It may convert one force to another force. - This is generally true as mechanical systems often convert forces between different parts.

Option C: It may convert one speed to another speed. - Mechanical systems often include gears and linkages that change speeds.

Option D: It may convert linear motion into rotational motion. - **This is explicitly shown in the video.**

Option E: It may convert gravitational energy to kinetic energy. - This transformation could be inferred based on the interactions of the parts, but it is not the main focus of the video.

After careful analysis, the final answer is: {A}.

**Error Category: Logical Reasoning Error**

**Error Reason:** There is self-contradiction in the entire reasoning process. The model believes that kinetic energy can be converted into elastic potential energy, and that gravitational potential energy can be converted into kinetic energy. According to this logic, gravitational potential energy can also be converted into elastic potential energy, which aligns with the statement in option A. However, the model incorrectly assumes that there is no mechanical transformation mentioned in option A in the video, resulting in a contradictory reasoning process.

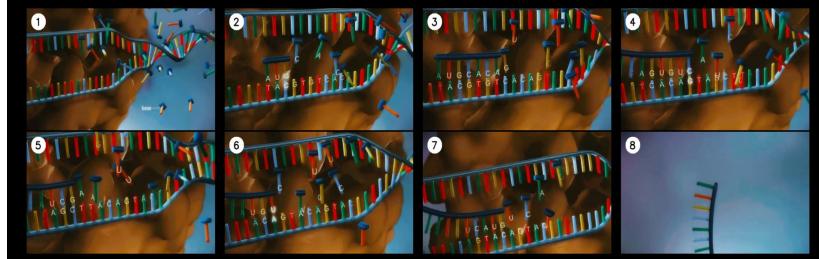
Figure 28: An error case of Mechanical Engineering.

## Clinical Medicine

**Question:** What could the brown stuff in the video be?

**Options:**

- (A) peptidyltransferase (B) RNA polymerase (C) DNA polymerase (D) Topoisomerase (E) Spliceosome complex



**Answer:** B

**Relevant Knowledge (Wikipedia page title):** RNA polymerase, Transcription (biology), Peptidyltransferase, DNA polymerase, Topoisomerase, Spliceosome

**Reasoning Rational:**

The video demonstrates the process of DNA unwinding and transcription of RNA under the action of enzymes, specifically highlighting the role of RNA polymerase.

### DeepSeek-VL2:

The brown structure in the video likely represents a biological process, as it is associated with DNA, RNA, and various molecular structures. The brown stuff could correspond to an enzyme that replicates DNA or RNA, or it could represent a complex that splices a strand of RNA. RNA polymerase synthesizes DNA strands. RNA polymerase synthesizes RNA from DNA, and peptidyltransferase catalyzes the formation of peptide bonds between two amino acids. Topoisomerase adjusts the tension or catenation in DNA by introducing or removing supercoils. Spliceosomes are complexes involved in pre-mRNA splicing.

Given the context, it seems most plausible that the brown stuff represents an enzyme involved in the replication of RNA ([RNA polymerase](#)) or an enzyme involved in cutting and splicing an mRNA strand ([spliceosome complex](#)). Therefore, we can eliminate options A, B, C, and D based on the actions shown in the video. ([Human Evaluator's Comment: Self-contradiction](#)) The most likely final answer is: E, Spliceosome complex.

### Error Category: Logical Reasoning Error

**Error Reason:** The model thinks that the brown stuff represents a RNA polymerase or a spliceosome complex. But when the model comes to the conclusion, it excludes the correct answer spliceosome complex without any further reason, which contradicts with its previous logical reasoning.

Figure 29: An error case of Clinical Medicine.