# Video-MMLU: A Massive Multi-Discipline Lecture Understanding Benchmark







Enxin Song <sup>1</sup>

Wenhao Chai <sup>2</sup>

<sup>1</sup>Zhejiang University

Weili Xu<sup>1,3</sup>

<sup>2</sup>Princeton University

Jianwen Xie <sup>4</sup>

<sup>3</sup>University of Illinois Urbana-Champaign

Yuxuan Liu<sup>3</sup>

Gaoang wang 1

<sup>4</sup>Lambda, Inc.



### Contributions

Video-MMLU pushes LMMs to the limits—can the model really understand real-world lectures?

- We build Video-MMLU, which requires strong reasoning capabilities and world knowledge compared to the previous benchmarks for video LMMs.
- We evaluate more than 90 proprietary models and open-source models of varying sizes on VideoMMLU. Our findings indicate that existing models generally perform poorly, with accuracy ranging from only 10% to 50%.
- We explore how the number of visual tokens and the base LLMs influence performance, offering insights into the interplay between multimodal perception and reasoning in lecture comprehension.

## **Benchmark Statistics**

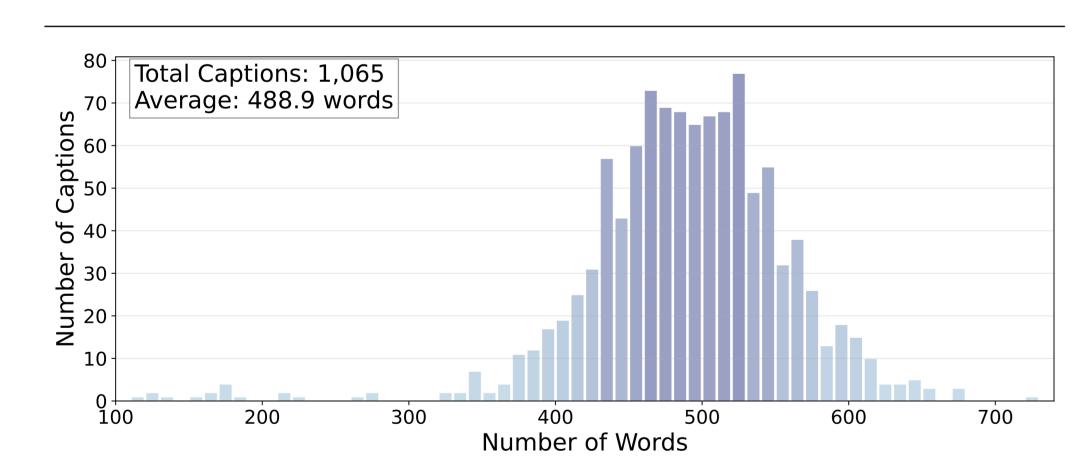


Figure 1. Video detailed captions length distribution.

8.0%	41.5%	19	9.1%	19.1%	12.4%
< 10frm	10frm~20frm	20frn	n~30frm	30frm~40frm	> 40frm
7.6%	30.6%	17.2%	26.8	3%	L7.8%
< 30s	30s~60s	60s~120s	120s~	180s	> 180s

Figure 2. Video length and keyframes number distribution.

#### Overview

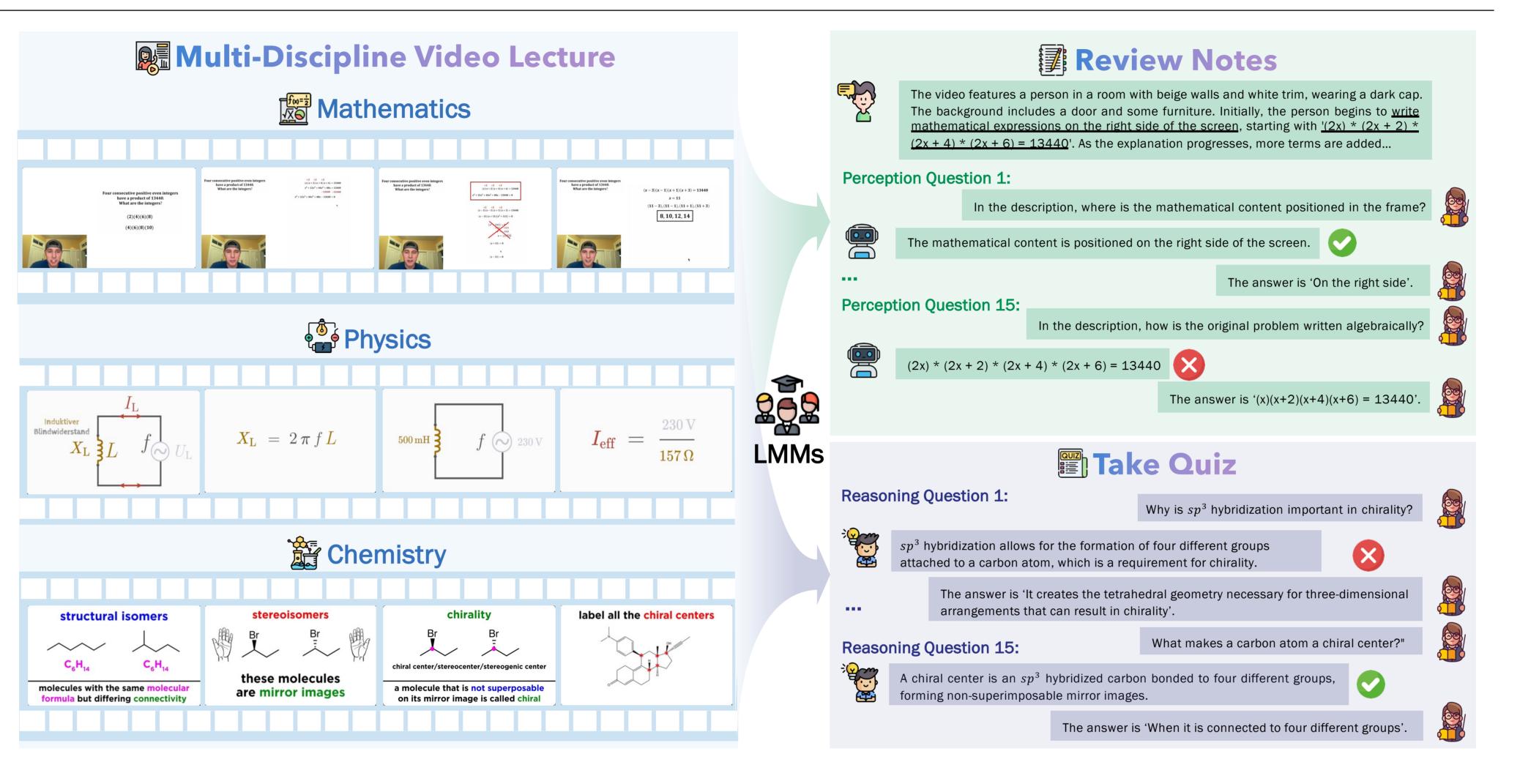


Figure 3. The benchmark includes multi-discipline lecture videos in mathematics, physics, and chemistry, featuring theorem demonstrations and problem-solving. Evaluation consists of: (1) **Review Notes**, where models generate detailed video captions to assess visual perception, and (2) **Take Quiz**, where models answer reasoning questions to test comprehension.

## Influence of model size

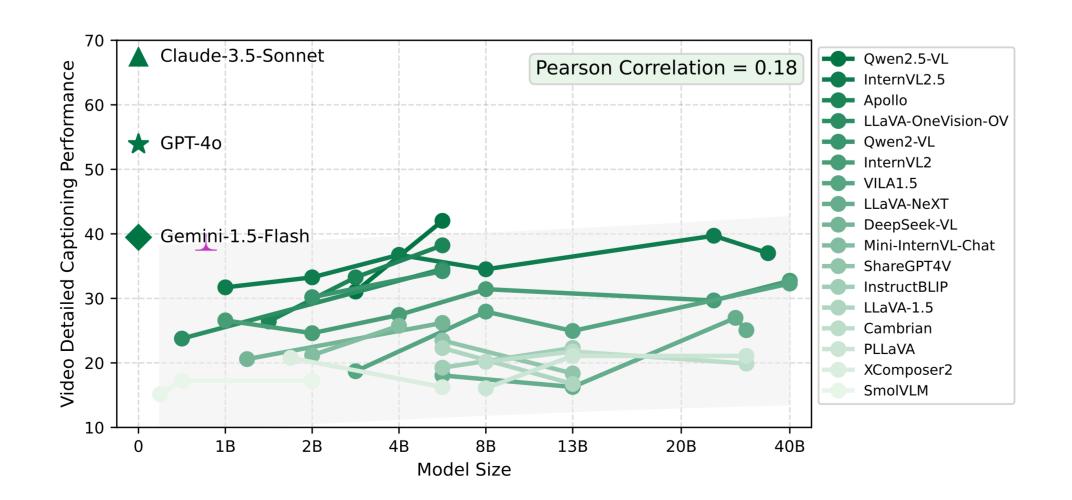


Figure 4. Relationship between model size and video captioning performance. The shaded region shows the confidence interval.

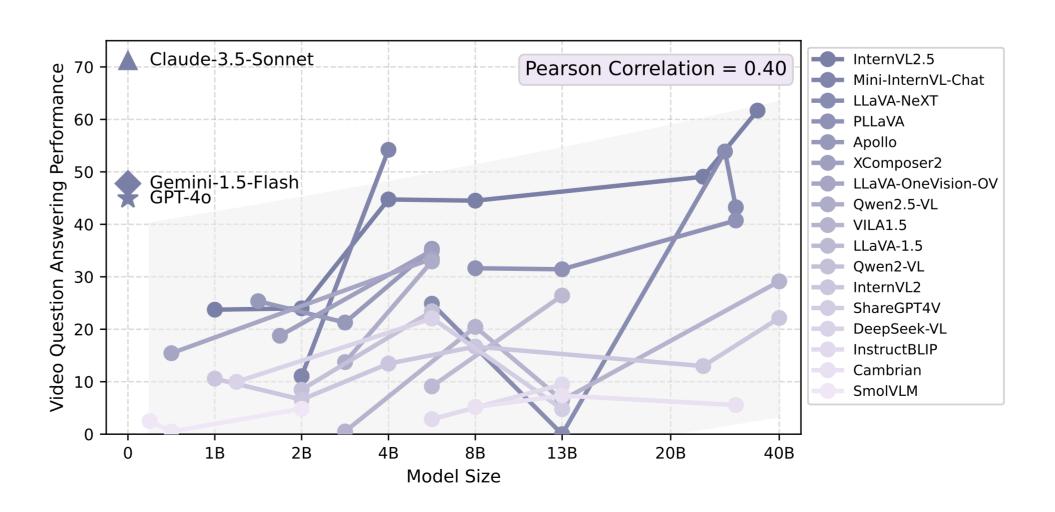


Figure 5. Relationship between model size and video QA performance. The shaded region shows the confidence interval.

#### Influence of LLM Architecture

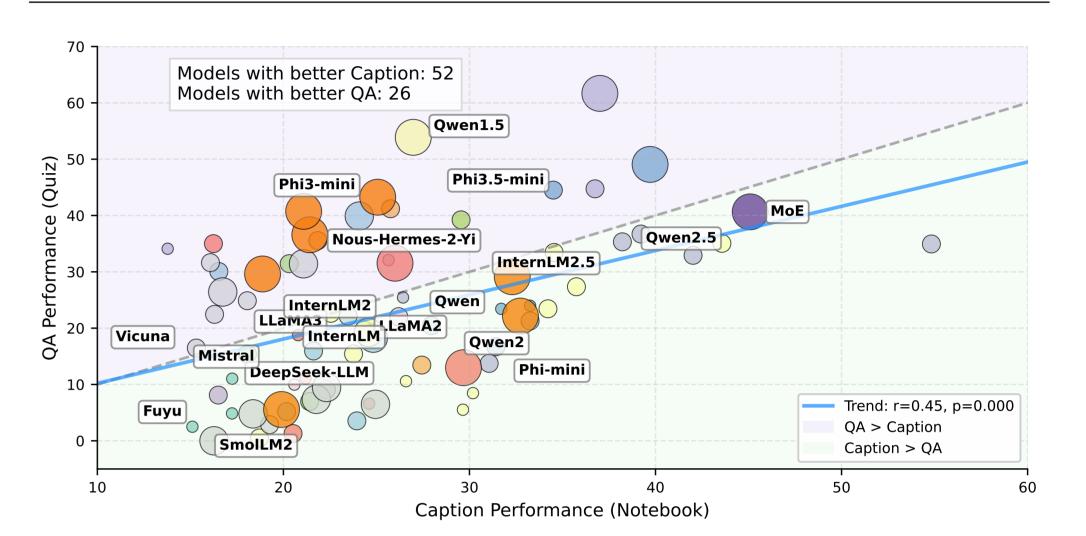


Figure 6. Relationship between captioning and QA across LLM.

## **Ability of Token Compression Models**

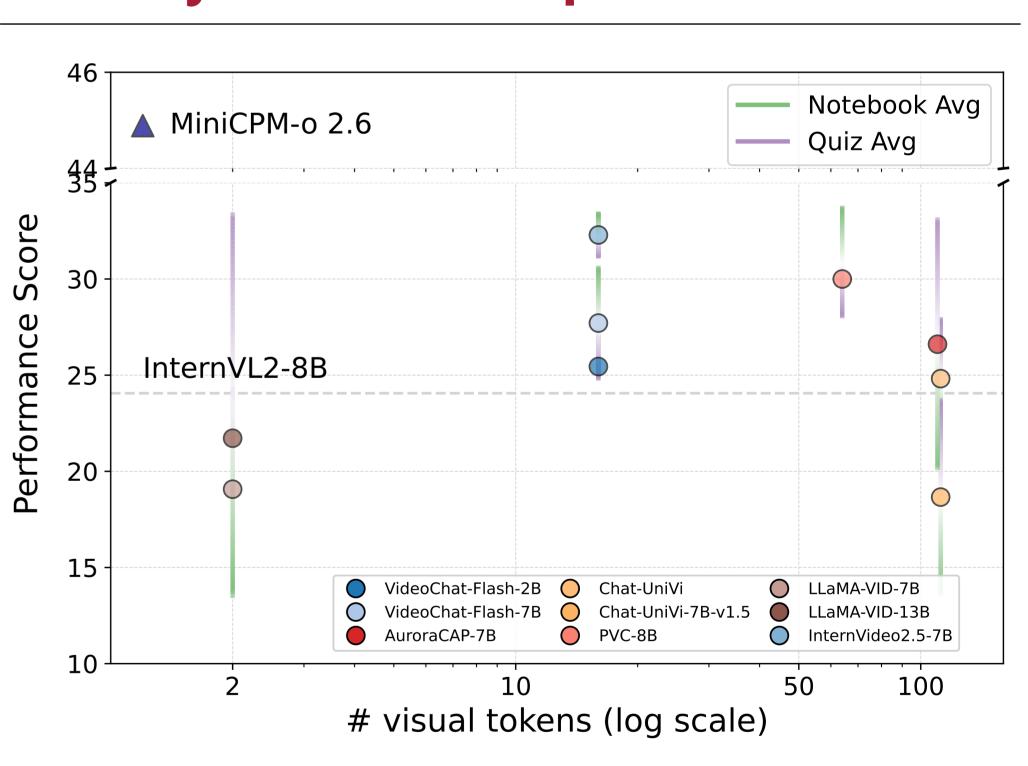


Figure 7. Impact of vision tokens number.

Our findings indicate that significant token reduction is feasible while maintaining or even surpassing the performance of the base model. Despite efficiency gains, a substantial gap remains between token-compressed models and the state-of-the-art, indicating challenges in preserving fine-grained details essential for complex lecture reasoning with existing token-compressed models.