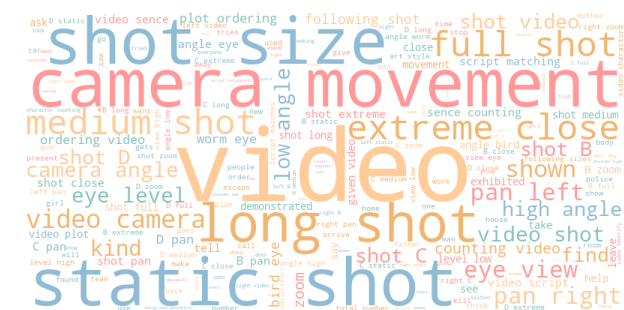
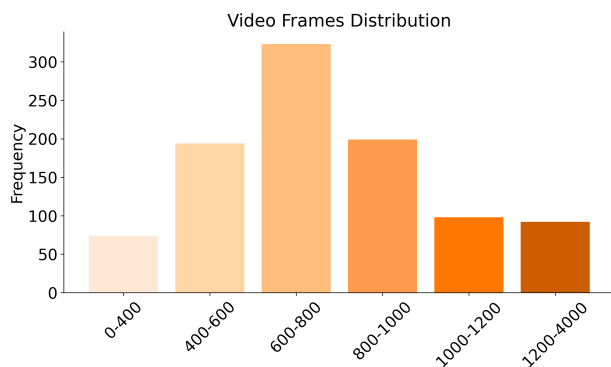


Supplementary Material

In this section, we show more statistics for VIDCOMPOSITION. Figure 6 presents a word cloud that highlights key



terms from the questions and options, demonstrating the diversity of video compositions included in our benchmark. Figure 7 illustrates the distribution of video frames across



different ranges, highlighting the diversity in video durations present in our dataset. Most videos are concentrated in the 600-800 frame range, with fewer videos having shorter or longer durations. This distribution reflects a balanced yet diverse set of videos, suitable for comprehensive benchmarking.

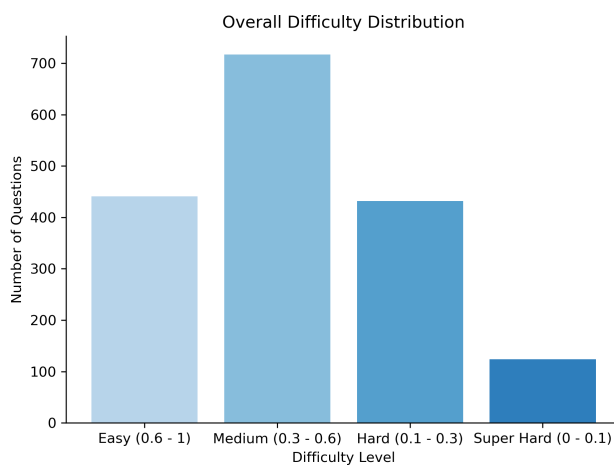


Figure 8. Distribution of questions across different difficulty levels, ranging from “Easy” (answered correctly by $>60\%$ of models,) to “Super Hard” (answered correctly by $<10\%$ of models).

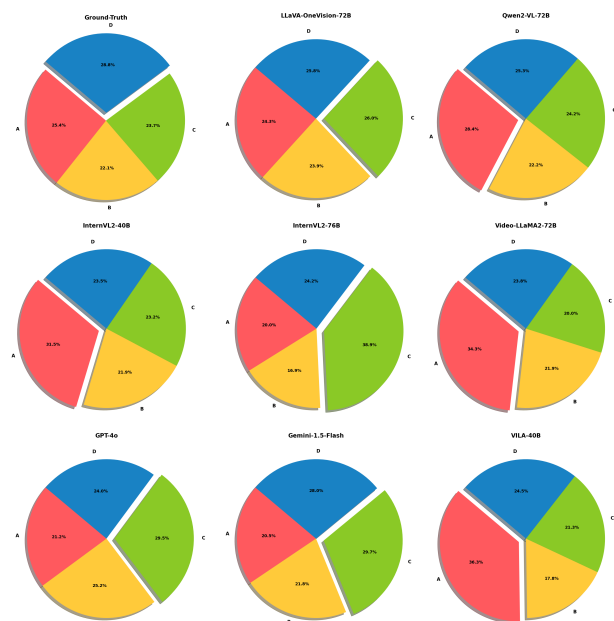


Figure 9. The distribution of answers for human-annotated questions is relatively balanced, as shown in the Ground-Truth pie chart. Predictions from several top models also exhibit relatively balanced distributions.

Figure 8 presents the distribution of questions across four difficulty levels: “Easy” (answered correctly by more than 60% of models), “Medium” (answered correctly by 30%-60% of models), “Hard” (answered correctly by 10%-30% of

models), and “Super Hard” (answered correctly by less than 10% of models). The pie charts in Figure 9 show that the answers for human-annotated questions are distributed fairly evenly across all options. Similarly, predictions from several top-performing models also demonstrate a comparable level of balance in their distributions. Comparing with the main results table, it can be observed that better model performance correlates with more evenly distributed predictions. For example, InternVL-40B outperforms InternVL-76B, as the pie chart reveals that InternVL-76B predictions are less evenly distributed and are biased toward option C compared to InternVL-40B. As shown in Figure 14, the tendency for imbalanced predictions is more pronounced in models with smaller sizes of LLM.

Figure 10 provides a detailed view of the distribution of questions across four difficulty levels (Easy, Medium, Hard, Super Hard) for each sub-task, reflecting the diverse challenges within the benchmark.

9. Tasks Definition in VIDCOMPOSITION

- **Camera Movement Perception (CamM-P)**: Identifying the types of camera movements shown in a video, such as panning, zooming, or tracking, which affect the visual flow and dynamics of the scene.
- **Shot Size Perception (SS-P)**: Recognizing the shot sizes, like close-up, medium shot, or full shot, which contribute to the viewer’s sense of intimacy or scope within the scene.
- **Camera Angle Perception (CamA-P)**: Identifying different camera angles used in the video, such as low angle, bird’s-eye view, or over-shoulder, which impact the perspective and interpretive context of the scene.
- **Emotion Perception (E-P)**: Detecting the emotions displayed by characters, such as fear, sadness, or happiness, which contribute to narrative understanding and character development.
- **Action Perception (A-P)**: Recognizing actions performed by characters, like driving or talking, to understand the physical activities and plot progression in the video.
- **Costume, Makeup, and Props Perception (CMP-P)**: Identifying elements of costume, makeup, and props used by characters, which provide contextual and stylistic cues about the setting, era, or genre.
- **Character Counting (Cha-C)**: Counting the number of characters appearing in the video, which gives an understanding of scene complexity and interaction density.
- **Script Matching (S-M)**: Identifying the narrative script or dialogue that corresponds with the visual content, facilitating alignment of visual and textual story elements.
- **Plot Ordering (P-O)**: Determining the chronological sequence of events depicted in the video, enabling coherent understanding of the storyline and causality.
- **Background Perception (B-P)**: Recognizing the type of background setting, such as a lakeside or grassland, which

anchors the scene’s location and environmental context.

- **Scene Counting (S-C)**: Counting the distinct scenes or settings within the video, indicating shifts in location or time that structure the narrative.
- **Lighting Perception (L-P)**: Identifying lighting conditions in the video, like high-key or low-key lighting, which affect the mood, visibility, and aesthetic of the scenes.
- **Art Style Perception (AS-P)**: Recognizing the art style of the video, such as Japanese cel anime or 3D CG animation, which contributes to the visual genre and artistic tone.
- **Cut Counting (Cut-C)**: Counting the number of cuts in the video, which reflects editing style and pacing, impacting the rhythm and viewer engagement.
- **Special Effect Perception (SE-P)**: Identifying special effects used in the video, like explosions or rain, which add dramatic or fantastical elements to enhance the visual experience.

10. Prompt Template

The prompt template provides explicit instructions for selecting the correct answer from multiple-choice options based on the provided video. It is carefully crafted to minimize ambiguity and guide the model’s reasoning process. By embedding a professional perspective (*e.g.*, “like a director and cinematographer”), the prompt attempts to align the model’s decision-making with human-like attention to detail. Additionally, the rigid structure ensures the compatibility of the generated responses, with only “A,” “B,” “C,” or “D” in the output.

Prompt Template for Model Prediction

Given a video, a multiple-choice question, and several options, ensure you select the option that correctly answers the question based on the provided video. Please consider comprehensively and meticulously like a professional director and cinematographer. Answer with the option’s letter from the given choices directly, and don’t contain any other contents!

{question}
{options}

In qualitative analysis, we also ask models to explain their answers. At this time, we replace “Answer with the option’s letter from the given choices directly, and don’t contain any other contents!” with “The output should contain the option index and explain why you selected this as your answer.”

11. Annotation & Human Evaluation System

The annotation checker user interface is shown as Figure 11. The system is designed to ensure high-quality annotations in the benchmark. The user interface displays the video,

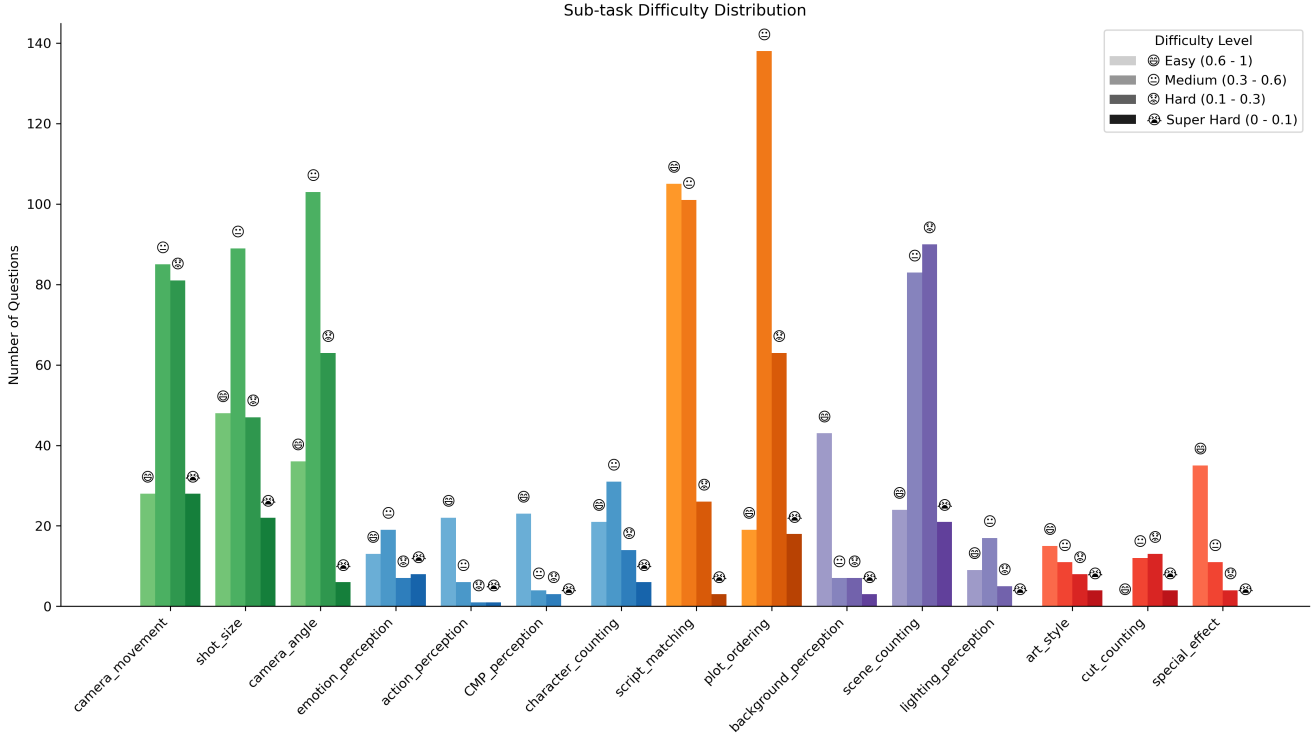


Figure 10. The difficulty distribution across various tasks in the benchmark. The bars represent the number of questions categorized into different difficulty levels (Easy, Medium, Hard, Super Hard) for each sub-task, highlighting variations in difficulty distribution across tasks.

the associated question, and multiple-choice options. It also includes a feedback section at the bottom, where reviewers can provide corrections or specify errors in the questions or options. Reviewers use this interface to assess the clarity and accuracy of annotations by attempting the questions themselves. If any issue is identified, such as unclear phrasing or incorrect answer options, they can use the feedback section to suggest improvements or note discrepancies. In addition to annotation refinement, this system is also used for human evaluation tasks, enabling consistent validation of both the dataset and the benchmark design. This iterative process ensures reliability, reduces errors, and supports the continual enhancement of question clarity and dataset quality.

12. More Results

In this section, we present the complete results of the diagnostic analysis on factors influencing the performance of MLLMs on VIDCOMPOSITION. Specifically, the impact of the number of frames is illustrated in Figure 12, the resolution of the visual encoder is detailed in Table 7, the size of the LLM is analyzed in Table 8, and the effect of training data volume is shown in Table 9. We also provide more visualization results in this section.

Watch the video and answer the following questions.

0:20 / 0:20

Show/Hide Link

Which of the following scale are not included in this film segment?

☐ A. close-up, long shot
 ☐ B. close-up, extreme close-up
 ☐ C. long shot, full shot
 ☐ D. extreme close-up, long shot
 ☒ If you believe none of these choices is correct, please enter your reason below.

Figure 11. The user interface for annotation checker. It can also be used for human evaluation.

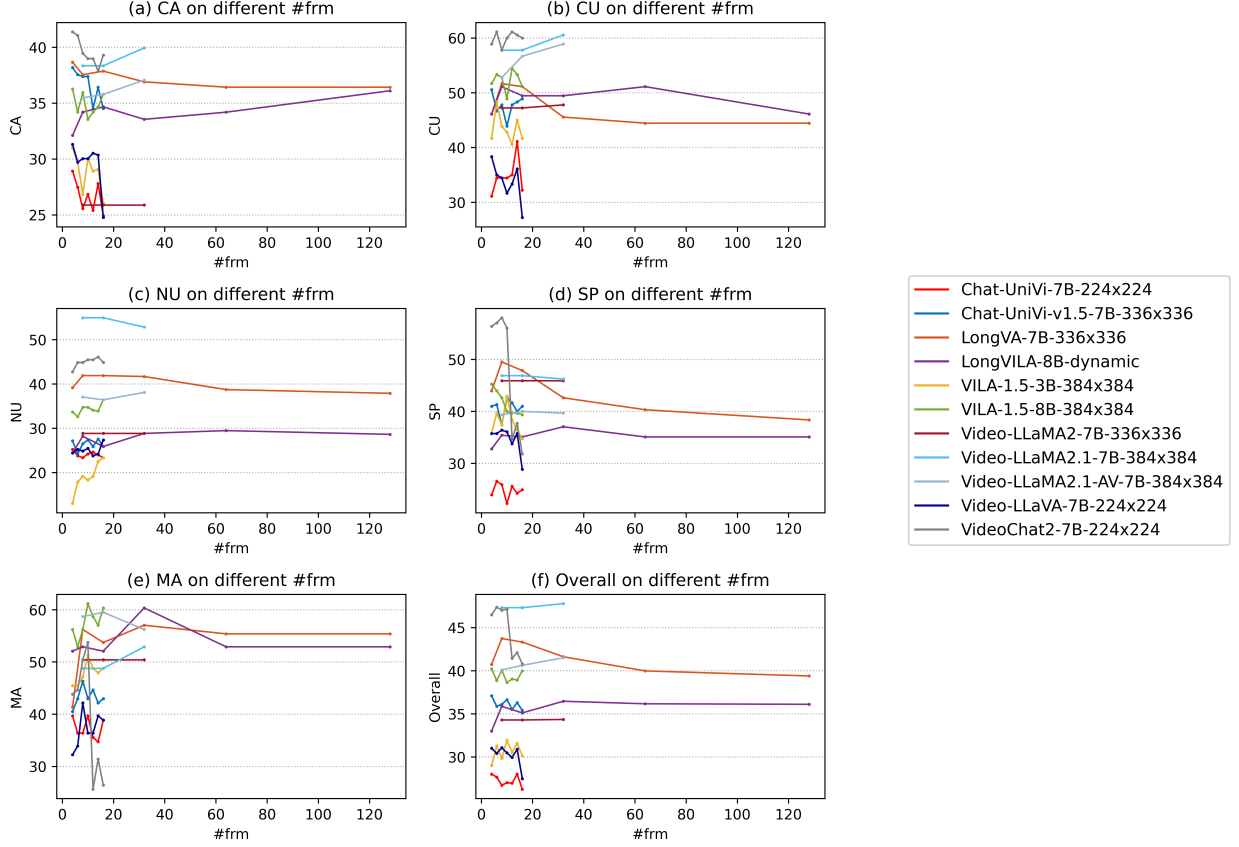


Figure 12. Full frame analysis. Performance analysis of models across different numbers of input frames. The results show no clear trends, with performance either remaining stable or fluctuating randomly as the number of frames increases.

Table 7. Full resolution analysis.

| Models | #frm | LLM size | Res. | CA | CU | NU | SP | MA | Overall |
|---------------------------------------|------|----------|------|--------------------------------|--------------------------------|--------------------------------|--------------------------------|--------------------------------|--------------------------------|
| Chat-UniVi; Video-LLaVA; VideoChat2 | 4 | 7B | 224 | 33.87 | 42.78 | 30.81 | 38.69 | 38.57 | 35.17 |
| Chat-UniVi-v1.5; LongVA | | | 336 | 38.42 ^{+4.55} | 48.33 ^{+5.55} | 33.16 ^{+2.35} | 42.46 ^{+3.77} | 40.91 ^{+2.34} | 38.92 ^{+3.75} |
| Chat-UniVi; Video-LLaVA; VideoChat2 | 6 | 7B | 224 | 32.75 | 43.52 | 31.3 | 39.78 | 38.29 | 35.15 |
| Chat-UniVi-v1.5 | | | 336 | 37.54 ^{+4.79} | 46.67 ^{+3.15} | 24.21 ^{-7.09} | 41.31 ^{+1.53} | 42.98 ^{+4.69} | 35.87 ^{+0.72} |
| Chat-UniVi; Video-LLaVA; VideoChat2 | 8 | 7B | 224 | 31.68 | 42.22 | 31.02 | 40.11 | 42.98 | 34.94 |
| Chat-UniVi-v1.5; LongVA; Video-LLaMA2 | | | 336 | 33.6 ^{+1.92} | 48.89 ^{+6.67} | 32.42 ^{+1.4} | 44.26 ^{+4.15} | 50.96 ^{+7.98} | 38.04 ^{+3.1} |
| Video-LLaMA2.1; Video-LLaMA2.1-AV | | | 384 | 36.9 ^{+3.3} | 55.28 ^{+6.39} | 46.0 ^{+13.58} | 43.11 ^{-1.15} | 53.72 ^{+2.76} | 43.7 ^{+5.66} |
| Chat-UniVi; Video-LLaVA; VideoChat2 | 10 | 7B | 224 | 31.95 | 42.04 | 31.72 | 38.14 | 43.25 | 34.88 |
| Chat-UniVi-v1.5 | | | 336 | 37.38 ^{+5.43} | 43.89 ^{+1.85} | 27.37 ^{-4.35} | 42.62 ^{+4.48} | 42.98 ^{-0.27} | 36.64 ^{+1.76} |
| Chat-UniVi; Video-LLaVA; VideoChat2 | 12 | 7B | 224 | 31.63 | 43.15 | 31.3 | 31.37 | 32.51 | 32.79 |
| Chat-UniVi-v1.5 | | | 336 | 34.5 ^{+2.87} | 47.78 ^{+4.63} | 25.89 ^{-5.41} | 41.64 ^{+10.27} | 44.63 ^{+12.12} | 35.52 ^{+2.73} |
| Chat-UniVi; Video-LLaVA; VideoChat2 | 14 | 7B | 224 | 32.0 | 45.93 | 31.44 | 32.57 | 35.26 | 33.67 |
| Chat-UniVi-v1.5 | | | 336 | 36.42 ^{+4.42} | 48.33 ^{+2.4} | 27.58 ^{-3.86} | 40.0 ^{+7.43} | 42.15 ^{+6.89} | 36.28 ^{+2.61} |
| Chat-UniVi; Video-LLaVA; VideoChat2 | 16 | 7B | 224 | 29.66 | 39.81 | 31.86 | 28.52 | 34.71 | 31.52 |
| Chat-UniVi-v1.5; LongVA; Video-LLaMA2 | | | 336 | 32.75 ^{+3.09} | 49.07 ^{+9.26} | 32.21 ^{+0.35} | 44.92 ^{+16.4} | 49.04 ^{+14.33} | 37.67 ^{+6.15} |
| Video-LLaMA2.1; Video-LLaMA2.1-AV | | | 384 | 37.06 ^{+4.31} | 57.22 ^{+8.15} | 45.68 ^{+13.47} | 43.44 ^{-1.48} | 54.13 ^{+5.09} | 43.96 ^{+6.29} |
| VILA-1.5 | 16 | 8B | 384 | 35.78 | 51.11 | 36.42 | 39.34 | 60.33 | 39.98 |
| InternVL2 | | | 448 | 57.03 ^{+21.25} | 62.78 ^{+11.67} | 53.68 ^{+17.26} | 45.57 ^{+6.23} | 56.2 ^{-4.13} | 54.63 ^{+14.65} |
| LongVA; Video-LLaMA2 | 32 | 7B | 336 | 31.39 | 46.67 | 35.26 | 44.26 | 53.72 | 37.98 |
| Video-LLaMA2.1; Video-LLaMA2.1-AV | | | 384 | 38.5 ^{+7.11} | 59.72 ^{+13.05} | 45.47 ^{+10.21} | 42.95 ^{-1.31} | 54.55 ^{+0.83} | 44.64 ^{+6.66} |

Human
 ① LLaVA-OneVision-72B
 InternVL2-40B
 Gemini-1.5-Flash
 GPT-4o

Full Shot
 Full Shot
 Medium Shot
 Close-up
 Close-up

Which of the following shot sizes are exhibited in this video?

A. Full shot, medium shot B. Long shot, close-up C. Extreme close-up, full shot D. Extreme close-up, long shot

① **B(long shot, close-up).** The video shows a **long shot** of the person walking down the alley and a **close-up** of the person standing in the classroom.

I chose **(B) long shot, close-up**, because the video shows a man walking down a street in a **long shot**, and then a **close-up** of a woman in a classroom. The **long shot** captures the man's entire body as he walks, while the **close-up** focuses on the woman's face and upper body.

The answer is **(D) extreme close-up, long shot**. The first shot of the video is a **long shot** showing the man walking down the street. Later, the video cuts to an **extreme close-up** on the woman's face as she speaks. **This demonstrates both long shots and extreme close-ups.**

B. The initial images show a wide view of an outdoor scene, characteristic of a **long shot** where the human figure is relatively small and the environment is prominent. Later images are tightly framed on the subject's face, typical of a **close-up**, which focuses on facial expressions and details.

Figure 13. Example of a multiple-choice question assessing **SS-P** in a video. The figure compares predictions from human annotators and various models, highlighting their selected answers and reasoning processes. **Correct answers** and explanations are provided for reference.

Table 8. Full LLM size analysis.

| Model | Res. | #frm | LLM size | CA | CU | NU | SP | MA | Overall |
|--------------|---------|-------|----------|--------------------------------|--------------------------------|--------------------------------|--------------------------------|--------------------------------|--------------------------------|
| InternVL2 | Dynamic | 16 | 0.5B | 24.12 | 28.33 | 26.53 | 29.84 | 28.93 | 26.61 |
| | | | 1.8B | 24.28 ^{+0.16} | 54.44 ^{+26.11} | 35.16 ^{+8.63} | 47.54 ^{+17.7} | 53.72 ^{+24.79} | 36.75 ^{+10.14} |
| | | | 3.8B | 32.11 ^{+7.83} | 51.67 ^{-2.77} | 44.0 ^{+8.84} | 48.85 ^{+1.31} | 48.76 ^{-4.96} | 41.68 ^{+4.93} |
| | | | 8B | 57.03 ^{+24.92} | 62.78 ^{+11.11} | 53.68 ^{+9.68} | 45.57 ^{-3.28} | 56.2 ^{+7.44} | 54.63 ^{+12.95} |
| | | | 20B | 40.1 ^{-16.93} | 63.89 ^{+1.11} | 51.16 ^{-2.52} | 38.36 ^{-7.21} | 54.55 ^{-1.65} | 46.42 ^{-8.21} |
| | | | 34B | 54.95 ^{+14.85} | 69.44 ^{+5.55} | 65.47 ^{+14.31} | 56.07 ^{+17.71} | 70.25 ^{+15.7} | 60.73 ^{+14.31} |
| | | | 70B | 51.92 ^{-3.03} | 72.78 ^{+3.34} | 64.84 ^{-0.63} | 52.79 ^{-3.28} | 63.64 ^{-6.61} | 58.73 ^{-2.0} |
| Qwen2-VL | Dynamic | 2 fps | 2B | 25.08 | 47.22 | 37.05 | 47.54 | 44.63 | 36.17 |
| | | | 7B | 34.35 ^{+9.27} | 56.67 ^{+9.45} | 61.05 ^{+24.0} | 57.38 ^{+9.84} | 48.76 ^{+4.13} | 49.3 ^{+13.13} |
| | | | 72B | 50.48 ^{+16.13} | 60.0 ^{+3.33} | 71.16 ^{+10.11} | 60.0 ^{+2.62} | 46.28 ^{-2.48} | 58.68 ^{+9.38} |
| VILA-1.5 | 384 | 4 | 3B | 30.99 | 41.67 | 13.05 | 35.74 | 45.45 | 29.02 |
| | | | 8B | 36.26 ^{+5.27} | 51.67 ^{+10.0} | 33.68 ^{+20.63} | 45.25 ^{+9.51} | 56.2 ^{+10.75} | 40.21 ^{+11.19} |
| | | 6 | 3B | 29.71 | 48.33 | 17.89 | 39.67 | 45.45 | 31.3 |
| | | | 8B | 34.19 ^{+4.48} | 53.33 ^{+5.0} | 32.63 ^{+14.74} | 43.93 ^{+4.26} | 52.89 ^{+7.44} | 38.86 ^{+7.56} |
| | | 8 | 3B | 26.84 | 43.89 | 19.16 | 37.38 | 47.11 | 29.84 |
| | | | 8B | 35.94 ^{+9.1} | 52.78 ^{+8.89} | 34.74 ^{+15.58} | 42.62 ^{+5.24} | 56.2 ^{+9.09} | 40.04 ^{+10.2} |
| | | 10 | 3B | 30.03 | 42.78 | 18.32 | 42.95 | 51.24 | 31.95 |
| | | | 8B | 33.55 ^{+3.52} | 48.89 ^{+6.11} | 34.74 ^{+16.42} | 40.0 ^{-2.95} | 61.16 ^{+9.92} | 38.63 ^{+6.68} |
| | | 12 | 3B | 28.91 | 40.56 | 19.16 | 38.36 | 48.76 | 30.54 |
| | | | 8B | 34.19 ^{+5.28} | 54.44 ^{+13.88} | 34.11 ^{+14.95} | 39.67 ^{+1.31} | 58.68 ^{+9.92} | 39.04 ^{+8.5} |
| | | 14 | 3B | 29.07 | 45.0 | 22.53 | 36.39 | 47.93 | 31.59 |
| | | | 8B | 34.66 ^{+5.59} | 53.33 ^{+8.33} | 33.89 ^{+11.36} | 39.67 ^{+3.28} | 57.02 ^{+9.09} | 38.92 ^{+7.33} |
| | | 16 | 3B | 26.04 | 41.67 | 23.37 | 34.75 | 48.76 | 30.13 |
| | | | 8B | 35.78 ^{+9.74} | 51.11 ^{+9.44} | 36.42 ^{+13.05} | 39.34 ^{+4.59} | 60.33 ^{+11.57} | 39.98 ^{+9.85} |
| Video-LLaMA2 | 336 | 32 | 7B | 25.88 | 47.78 | 28.84 | 45.9 | 50.41 | 34.35 |
| | | | 72B | 54.15 ^{+28.27} | 71.67 ^{+23.89} | 65.68 ^{+36.84} | 48.52 ^{+2.62} | 59.5 ^{+9.09} | 58.62 ^{+24.27} |

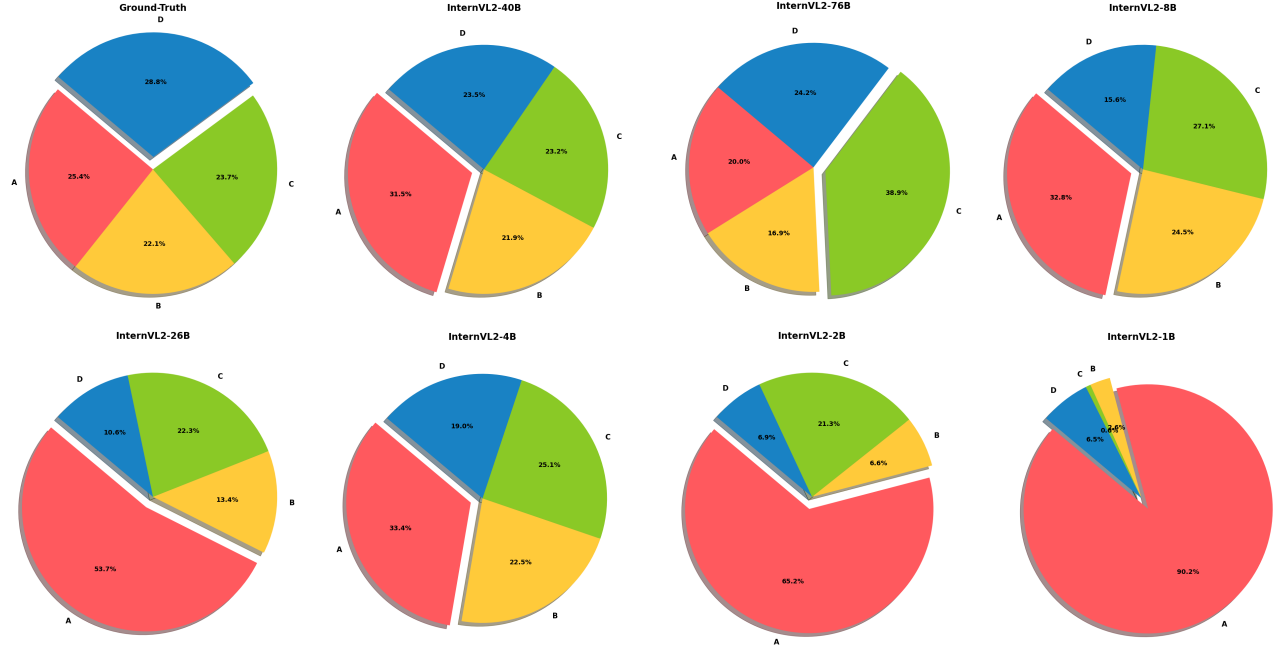


Figure 14. The pie charts show the answer distributions of InternVL models with varying LLM sizes. Larger models usually predict more balanced distributions, while smaller models, like InternVL-1B, exhibit strong biases, particularly toward option A.

Table 9. Full Data volume analysis.

| Model | #frm | Res. | LLM size | Data volume | CA | CU | NU | SP | MA | Overall |
|---|------|------|----------------|-------------------------|--|--|--|--|--|---|
| Chat-UniVi VideoChat2 | 4 | 224 | 7B | 0.65M 2M | 28.91 41.37 _{+12.46} | 31.11 58.89 _{+27.78} | 25.26 42.74 _{+17.48} | 23.93 56.39 _{+32.46} | 39.67 43.8 _{+4.13} | 28.02 46.48 _{+18.46} |
| Chat-UniVi-v1.5 LongVA | 4 | 336 | 7B | 1.27M 1.32M | 38.18 38.66 _{+0.48} | 50.56 46.11 _{-4.45} | 27.16 39.16 _{+12.0} | 40.98 43.93 _{+2.95} | 40.5 41.32 _{+0.82} | 37.1 40.74 _{+3.64} |
| Chat-UniVi VideoChat2 | 6 | 224 | 7B | 0.65M 2M | 27.48 41.05 _{+13.57} | 34.44 61.11 _{+26.67} | 23.79 44.84 _{+21.05} | 26.56 57.05 _{+30.49} | 36.36 44.63 _{+8.27} | 27.67 47.36 _{+19.69} |
| Chat-UniVi VideoChat2 | 8 | 224 | 7B | 0.65M 2M | 25.56 39.46 _{+13.9} | 34.44 57.78 _{+23.34} | 23.37 44.84 _{+21.47} | 25.9 58.03 _{+32.13} | 36.36 50.41 _{+14.05} | 26.73 47.01 _{+20.28} |
| Chat-UniVi-v1.5 LongVA | 8 | 336 | 7B | 1.27M 1.32M | 37.38 37.54 _{+0.16} | 47.78 51.67 _{+3.89} | 26.53 41.89 _{+15.36} | 37.38 49.51 _{+12.13} | 46.28 56.2 _{+9.92} | 36.11 43.73 _{+7.62} |
| VILA-1.5 Video-LLaMA2.1-AV Video-LLaMA2.1 | 8 | 384 | 8B 7B 7B | 1.21M 3.35M 3.35M | 35.94 35.46 _{-0.48} 38.34 _{+2.88} | 52.78 52.78 ₀ 57.78 _{+5.0} | 34.74 37.05 _{+3.31} 54.95 _{+17.9} | 42.62 39.34 _{-3.28} 46.89 _{+7.75} | 56.2 58.68 _{+2.48} 48.76 _{-9.92} | 40.04 40.09 _{+0.05} 47.3 _{+7.22} |
| Chat-UniVi VideoChat2 | 10 | 224 | 7B | 0.65M 2M | 26.84 38.98 _{+12.14} | 34.44 60.0 _{+25.56} | 24.21 45.47 _{+21.26} | 22.3 56.07 _{+33.77} | 39.67 53.72 _{+14.05} | 27.02 47.13 _{+20.11} |
| Chat-UniVi VideoChat2 | 12 | 224 | 7B | 0.65M 2M | 25.4 38.98 _{+13.58} | 35.0 61.11 _{+26.11} | 24.63 45.47 _{+20.84} | 25.57 34.75 _{+9.18} | 35.54 25.62 _{-9.92} | 26.96 41.44 _{+14.48} |
| Chat-UniVi VideoChat2 | 14 | 224 | 7B | 0.65M 2M | 27.8 37.86 _{+10.06} | 41.11 60.56 _{+19.45} | 24.0 46.11 _{+22.11} | 24.26 37.7 _{+13.44} | 34.71 31.4 _{-3.31} | 28.02 42.09 _{+14.07} |
| Chat-UniVi VideoChat2 | 16 | 224 | 7B | 0.65M 2M | 24.92 39.3 _{+14.38} | 32.22 60.0 _{+27.78} | 23.37 44.84 _{+21.47} | 24.92 31.8 _{+6.88} | 38.84 26.45 _{-12.39} | 26.26 40.8 _{+14.54} |
| Chat-UniVi-v1.5 LongVA | 16 | 336 | 7B | 1.27M 1.32M | 34.5 37.86 _{+3.36} | 48.89 51.11 _{+2.22} | 25.89 41.89 _{+16.0} | 40.98 47.87 _{+6.89} | 42.98 53.72 _{+10.74} | 35.4 43.32 _{+7.92} |
| VILA-1.5 Video-LLaMA2.1-AV Video-LLaMA2.1 | 16 | 384 | 8B 7B 7B | 1.21M 3.35M 3.35M | 35.78 35.78 ₀ 38.34 _{+2.56} | 51.11 56.67 _{+5.56} 57.78 _{+1.11} | 36.42 36.42 ₀ 54.95 _{+18.53} | 39.34 40.0 _{+0.66} 46.89 _{+6.89} | 60.33 59.5 _{-0.83} 48.76 _{-10.74} | 39.98 40.62 _{+0.64} 47.3 _{+6.68} |
| Kangaroo MiniCPM-V | 64 | 448 | 8B | 2.94M 8.32M | 31.79 38.18 _{+6.39} | 51.67 60.0 _{+8.33} | 29.05 40.84 _{+11.79} | 53.44 38.36 _{-15.08} | 33.06 55.37 _{+22.31} | 37.1 42.5 _{+5.4} |

Parse single letter choice

```
def extract_single_content(text):
    # If text is a list, convert it to a string by taking the first element
    if isinstance(text, list):
        if text: # Ensure the list is not empty
            text = text[0]
        else:
            return random.choice(['A', 'B', 'C', 'D']) # Return default if list is empty

    # Check if text is a valid string
    if not isinstance(text, str):
        return random.choice(['A', 'B', 'C', 'D'])

    # 1. Match patterns like (A)(B)(C)(D)
    match = re.search(r'\((A|B|C|D)\)', text)
    if match:
        return match.group(1)

    # 2. Match text starting with A, B, C, or D, followed by spaces or non-alphabetic characters
    match = re.match(r'^(A|B|C|D)[\s\W]*', text)
    if match:
        return match.group(1)

    # 3. Match standalone A, B, C, or D
    match = re.match(r'^\b[A-D]\b', text)
    if match:
        return match.group(1)

    # 4. Match patterns like (a), (b), (c), (d) and convert to uppercase
    match = re.search(r'\((a|b|c|d)\)', text)
    if match:
        return match.group(1).upper()

    # 5. Match patterns like A., B., C., or D.
    match = re.search(r'^\b(A|B|C|D)\.', text)
    if match:
        return match.group(1)

    # 6. If text contains a single letter, return it in uppercase
    letters = re.findall(r'[a-zA-Z]', text)
    if len(letters) == 1:
        return letters[0].upper()

    # Default return if no patterns match
    return random.choice(['A', 'B', 'C', 'D'])
```
