

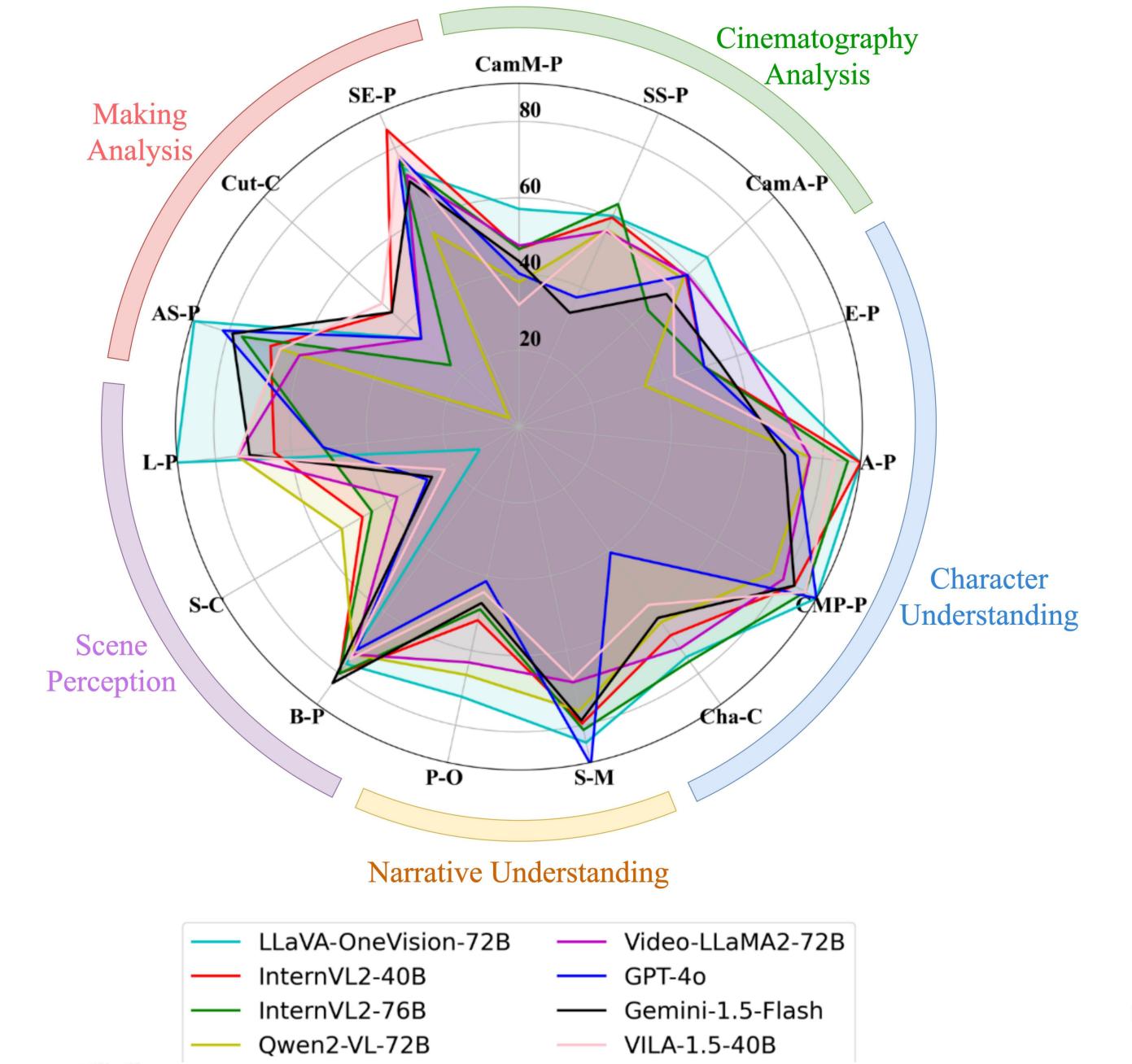
VidComposition: Can MLLMs Analyze Compositions in Compiled Video?

Yolo Y. Tang*, Junjia Guo*, Hang Hua, Susan Liang, Mingqian Feng, Xinyang Li, Rui Mao, Chao Huang, Jing Bi, Zeliang Zhang, Pooyan Fazli, Chenliang Xu

Introduction

- We introduce VidComposition, a novel and high-quality benchmark for evaluating fine-grained video composition understanding in MLLMs.
- We comprehensively evaluate 33 MLLMs for video understanding with VidComposition. The results show the challenging nature of VidComposition and the substantial gap between MLLMs' and humans' capabilities in video composition understanding.
- We analyze the critical factors that influence the performance of MLLMs systematically, providing potential directions for model improvement and future advancements.

Benchmark	I/V	#Data	#Task	Compositional QA	Compiled Videos	Fine-Grained Anno.
Winoground [49]	I	400	8	✓	-	X M
MME [11]	I	1.1k	14	X	-	✓ M
MMBench [34]	I	1.7k	20	X	-	✓ A+M
MMComposition [17]	I	4.3k	13	✓	-	✓ M
MSVD-QA [55]	V	504	5	X	X	X A
MSRVT-QA [55]	V	2.9k	5	X	X	X A
TGIF-QA [19]	V	9.6k	4	X	X	X A
TVQA [22]	V	2.2k	8	✓	✓	X A+M
ActivityNet-QA [59]	V	5.8k	4	X	X	X M
NExT-QA [54]	V	1k	8	X	X	X A
AutoEval-Video [6]	V	327	9	X	X	X A+M
Video-Bench [41]	V	5.9k	10	X	X	X A+M
LVBench [51]	V	500	6	X	X	X M
MVBench [26]	V	3.6k	20	X	X	✓ A+M
Movie-Chat-1k [45]	V	100	8	X	✓	✓ M
TempCompass [33]	V	410	4	X	X	✓ M
Video-MME [13]	V	900	12	X	X	✓ M
VidComposition	V	982	15	✓	✓	M



VidComposition

Cinematography Analysis

- What kind of movements of camera are shown in this video? A. zoom in, pan left B. pan right, pan down C. zoom out, pan left D. static shot, pan up
- ① Camera Movement Perception

Character Understanding

- Identify the emotion shown in the video. ④ Emotion Perception A. fear B. sad C. surprise D. happiness
- What actions can be seen in the video? A. driving a vehicle B. running C. talking to someone D. all of the above ⑤ Action Perception

Narrative Understanding

- Which script corresponds with this video? A. James looked angry at Anna in the submarine for ... B. A few days later, Rex's funeral was held ... C. Anna tried to open the safe and retrieve the nanomites bomb ... D. The Doctor orders his men to destroy the ice ...
- ⑧ Script Matching

Scene Perception

- What background is depicted in the video? A. lakeside B. grassland C. 3D CG Animation D. snow-covered landscape ⑩ Background Perception
- Can you identify the art style of this video? A. Japanese Cel Animation B. 3D Rendered 2D Look C. 3D CG Animation D. American Cel Animation ⑪ Art Style Perception

Making Analysis

- What's the total number of cuts in the given video? A. 9 B. 15 C. 3 D. 21 ⑫ Cut Counting
- How many distinct scenes are present in the video? A. 4 B. 10 C. 8 D. 12 ⑬ Scene Counting
- What is the lighting condition in the video? ⑫ Lighting Perception A. high-key lighting, natural lighting B. low-key lighting, artificial lighting C. natural lighting, low-key lighting D. all of the above
- What special effect is depicted in the video? ⑮ Special Effect Perception A. snow B. rain C. tornado D. explosion ⑯ Plot Ordering

Main Results and Analysis

Method	Cinematography Analysis			Character Understanding			Narrative Underst.			Scene Perception			Making Analysis			Overall
	CamM-P	SS-P	CamA-P	E-P	A-P	CMP-P	Cha-C	S-M	P-O	B-P	S-C	L-P	AS-P	Cut-C	SE-P	
Human	84.1	85.4	80.0	82.6	92.3	92.9	94.1	97.0	97.5	94.4	80.2	81.8	85.7	87.5	94.7	86.26
LLaVA-OneVision-72B [24]	57.1	60.5	66.3	63.3	90.0	90.0	74.6	84.7	72.4	76.9	12.0	90.3	89.5	34.5	74.1	63.31
InternVL2-40B [7]	46.6	60.0	58.9	51.0	90.0	83.3	67.6	79.6	51.9	80.0	47.4	64.5	68.4	44.8	85.2	60.73
InternVL2-76B [7]	46.6	63.9	45.5	51.0	86.7	86.7	76.1	81.3	49.0	80.0	44.5	51.6	76.3	24.1	75.9	58.73
Qwen2-VL-72B [50]	37.9	56.6	57.9	34.7	76.7	76.7	63.4	76.2	66.5	73.8	53.6	74.2	65.8	3.4	55.6	58.68
Video-LLaMA2-72B [9]	47.5	56.1	59.4	63.3	76.7	80.0	71.8	68.5	63.2	73.8	36.8	74.2	60.5	34.5	72.2	58.62
InternVL2-8B [7]	55.3	56.6	59.4	44.9	80.0	83.3	59.2	67.2	40.6	78.5	32.5	64.5	52.6	31.0	72.2	54.63
GPT-4o [1]	40.2	37.1	59.4	51.0	73.3	90.0	40.8	90.6	41.4	72.3	27.8	51.6	81.6	34.5	77.8	52.93
Gemini-1.5-Flash [44]	43.4	32.7	52.0	55.1	70.0	48.4	62.0	78.7	47.3	83.1	26.3	71.0	78.9	44.8	70.4	52.40
VILA-1.5-40B [29]	32.0	56.6	54.5	42.9	83.3	86.7	57.7	67.7	44.4	75.4	22.5	74.2	65.8	48.3	77.8	51.23
GPT-4 mini [1]	33.8	49.8	50.5	49.0	80.0	90.0	31.0	79.6	41.4	66.2	26.8	61.3	76.3	20.7	79.6	50.23
Gemini-1.5-Pro [44]	33.8	51.7	51.0	73.3	80.0	70.4	47.7	36.4	73.8	37.3	71.0	84.2	75.9	49.36		
Qwen2-VL-7B [50]	20.1	46.8	37.1	38.8	70.0	76.7	54.9	73.2	49.4	72.3	52.2	61.3	42.1	17.2	70.4	49.30
Oryx-7B [35]	34.7	54.1	57.4	57.1	80.0	73.3	66.2	48.5	34.7	73.8	41.6	61.3	39.5	20.7	66.7	48.77
Gemini-1.5-Flash-8B [44]	43.4	45.9	56.9	36.7	70.0	76.7	35.2	69.8	36.0	73.8	26.8	74.4	71.1	24.1	64.8	48.59
Video-LLaMA2.1 [9]	44.3	35.6	39.6	51.0	76.7	83.3	50.7	60.9	45.2	75.4	35.4	58.1	36.8	20.7	81.5	47.77
Video-Chat2 [26]	24.2	58.0	42.1	44.9	66.7	83.3	60.6	62.6	27.6	73.8	55.0	35.5	50.0	10.3	59.3	47.36
InternVL2-26B [7]	33.3	47.8	39.6	55.1	76.7	83.3	56.3	68.9	33.9	76.9	25.4	45.2	34.2	41.4	75.9	46.42
LongVA [62]	24.7	41.0	48.0	40.8	70.0	73.3	42.3	51.9	32.2	72.3	42.6	48.4	52.6	34.5	70.4	43.73
MiniCPM-V2.6 [56]	28.3	43.4	43.6	53.1	73.3	80.0	50.7	59.1	23.0	75.4	22.0	71.0	57.9	20.7	72.2	42.49
InternVL2-4B [7]	27.4	42.9	26.2	32.7	66.7	73.3	49.3	60.4	28.0	78.5	41.6	35.5	44.7	10.3	72.2	41.68
Video-LLaMA2.1-AV [9]	27.4	45.9	38.6	55.1	73.3	76.7	47.9	46.4	30.1	81.5	25.8	45.2	34.2	34.5	83.3	41.50
VILA-1.5-8B [29]	31.5	40.0	37.6	51.0	63.3	66.7	40.8	40.9	26.8	70.8	37.8	41.9	60.5	44.8	59.3	40.21
GPT-4-turbo [1]	23.7	37.1	35.1	46.9	63.3	80.0	25.4	54.9	36.4	50.8	29.7	64.5	39.5	44.8	70.4	39.85
LongLVA [53]	28.3	37.1	27.2	24.5	60.0	56.7	54.9	48.1	32.6	61.5	38.3	41.9	26.3	24.1	66.7	38.45
Kangaroo [31]	29.2	42.0	24.3	30.6	56.7	66.7	57.7	31.5	26.8	67.7	47.8	61.3	21.1	6.9	55.6	37.10
InternVL2-2B [7]	23.7	24.4	24.8	36.7	76.7	63.3	53.5									