





FACE-HUMAN-BENCH

A Comprehensive Benchmark of Face and Human Understanding for Multi-modal Assistants

Lixiong Qin^{1,*}, Shilong Ou^{1,*}, Miaoxuan Zhang^{1,*}, Jiangning Wei^{1,*}, Yuhang Zhang^{1,*}, Xiaoshuai Song¹, Yuchen Liu¹, Mei Wang², Weiran Xu¹

¹Beijing University of Posts and Telecommunications

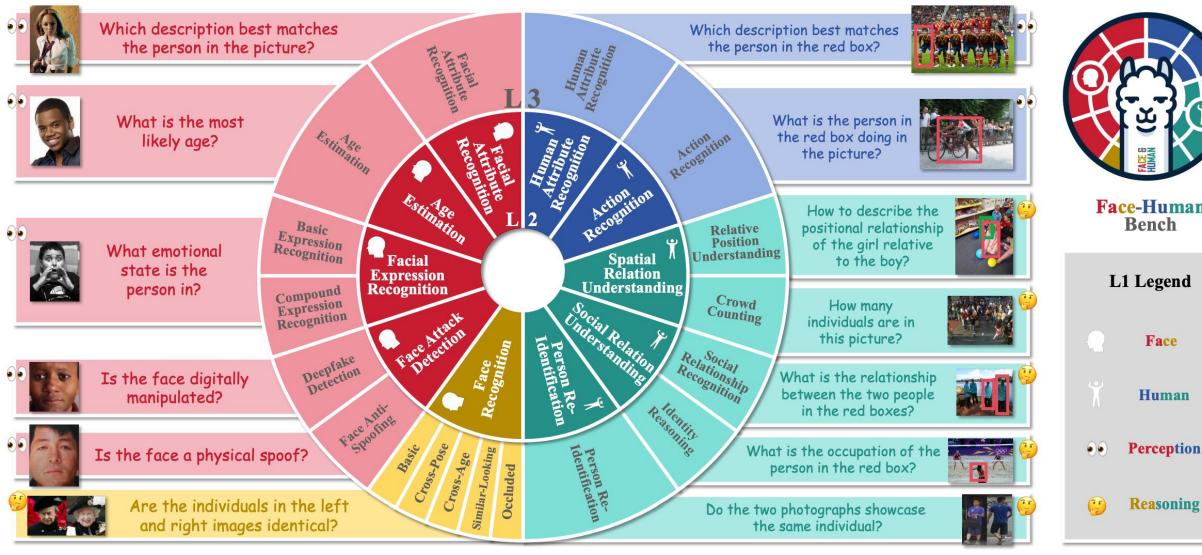
²Beijing Normal University

* Equal Contribution







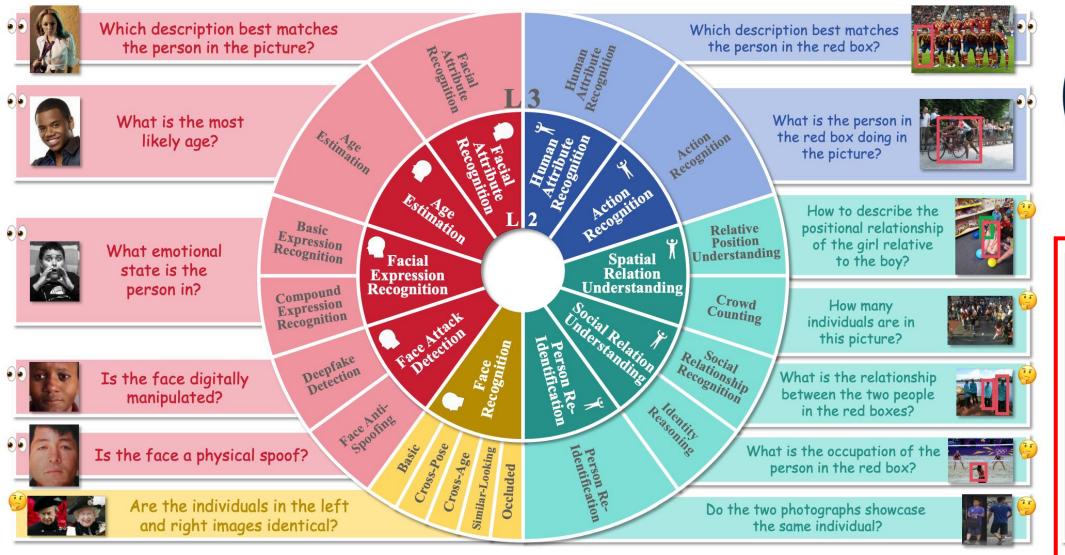


Human











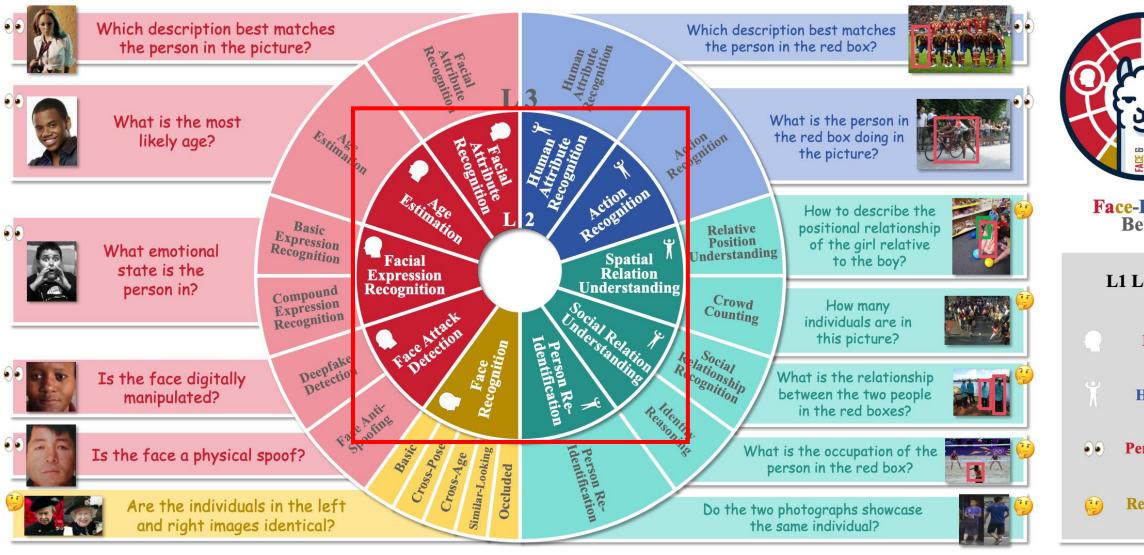
Face-Human Bench











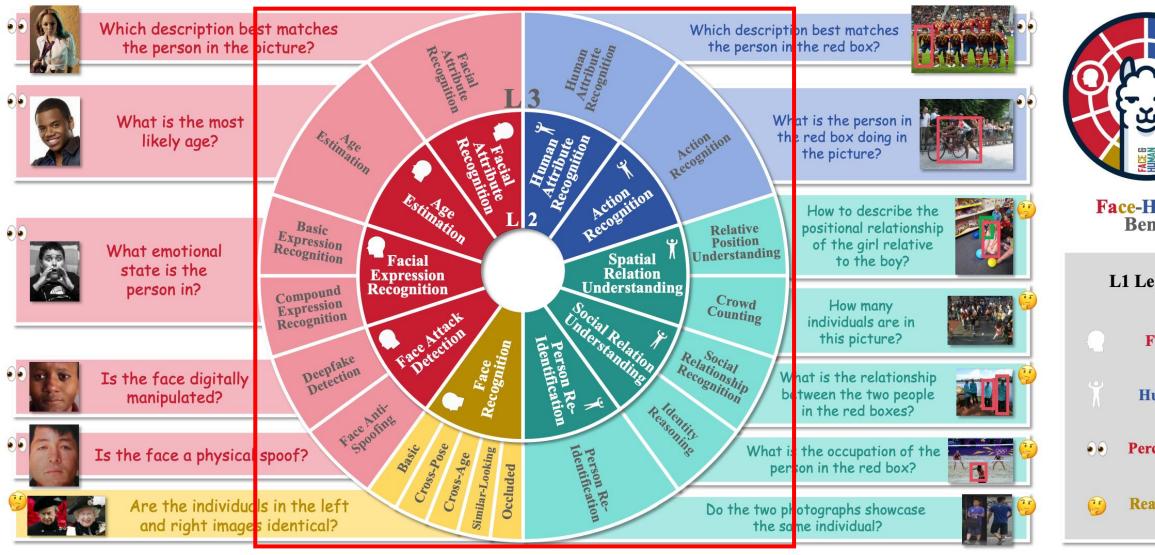
Face-Human Bench













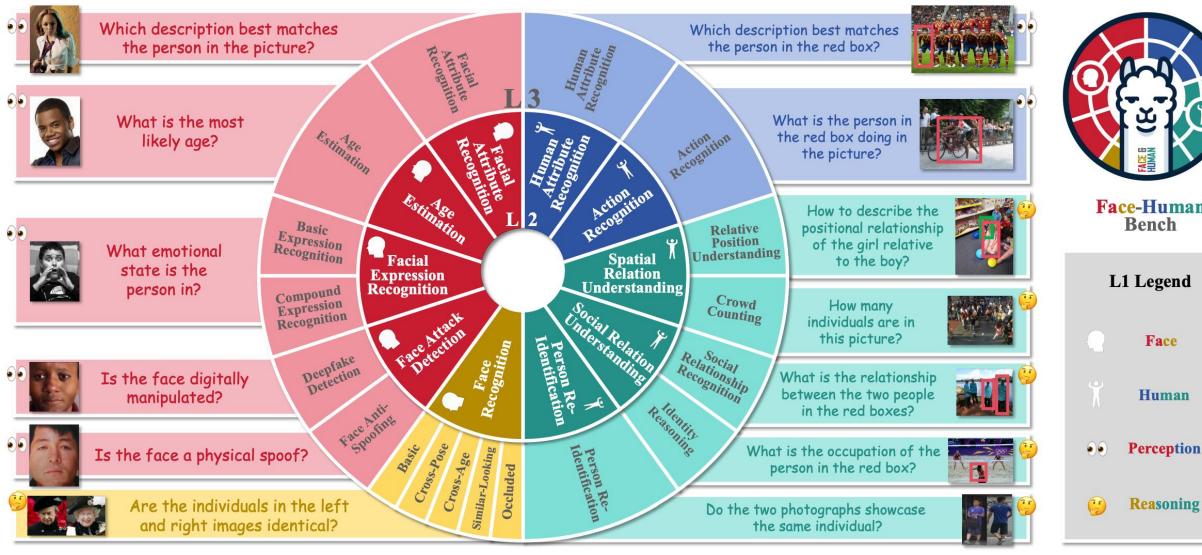
Face-Human Bench











Human



Our Research Questions





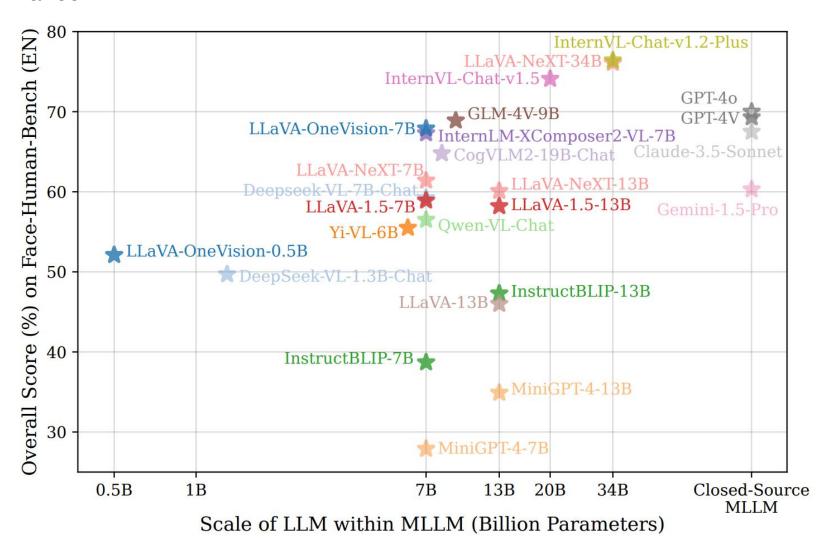
- Q1: How do existing MLLMs perform in face and human understanding?
 - 25 MLLMs' overall performance
 - Ability correlations across levels
 - Target relative position impacts
 - CoT prompting effects.

Q2: Which tasks do specialist models outperform MLLMs at?





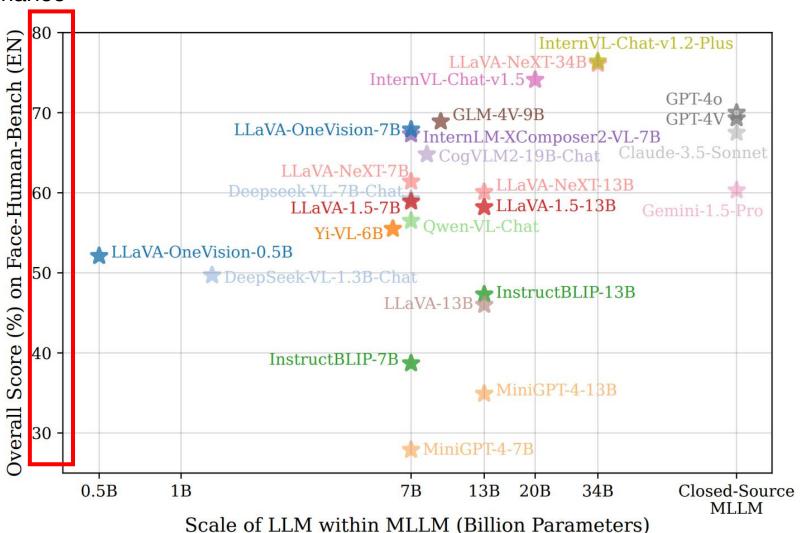








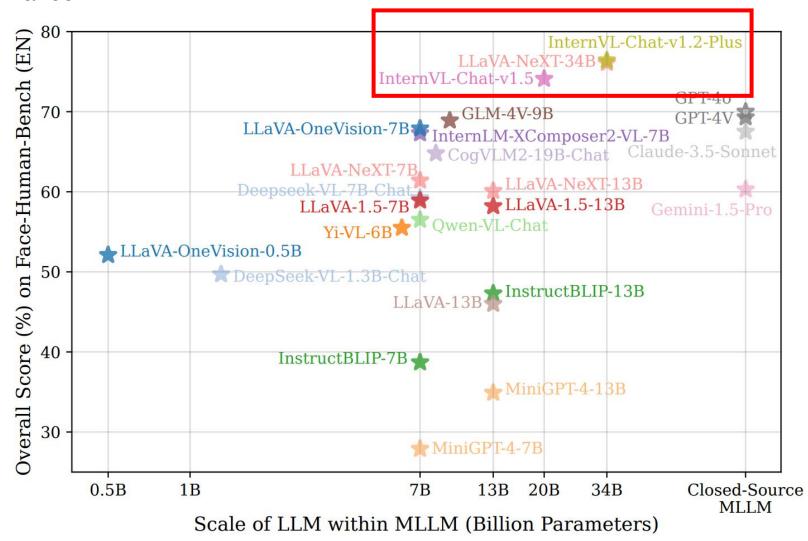








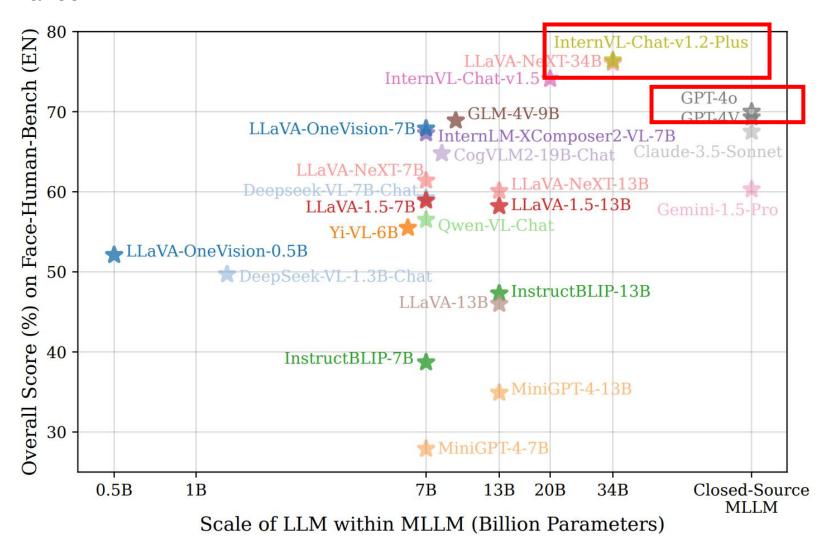








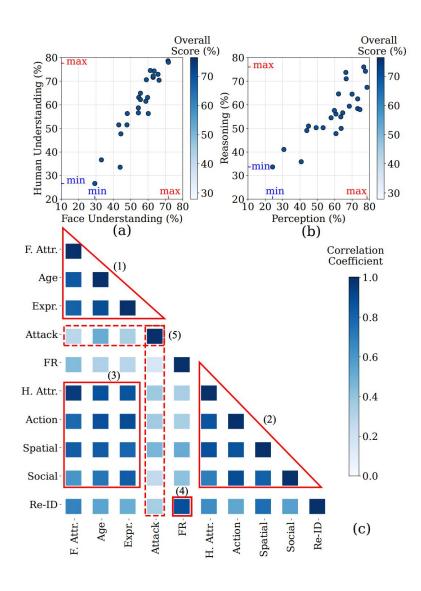








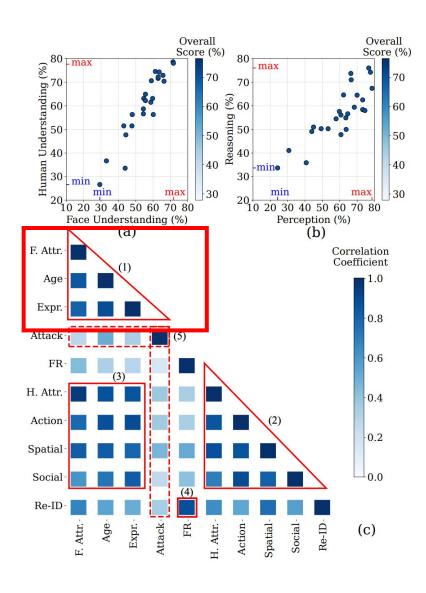








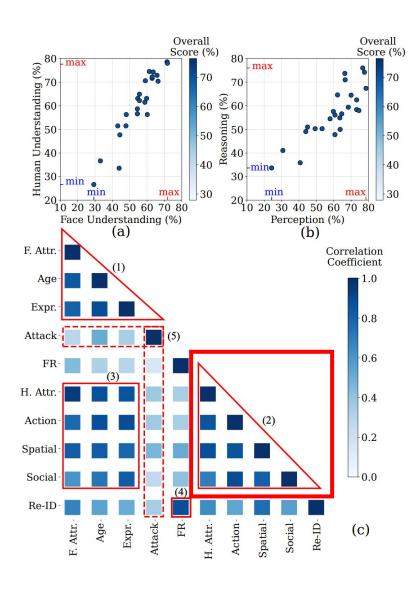








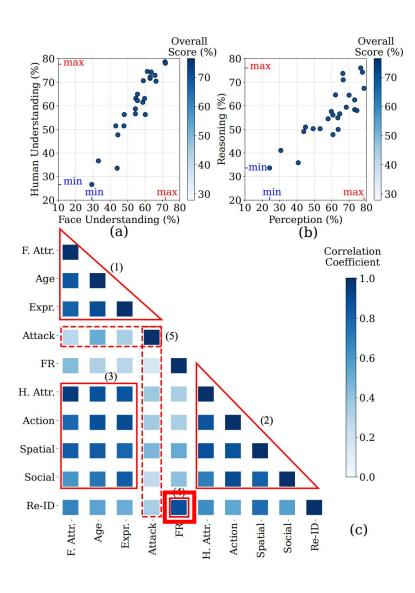










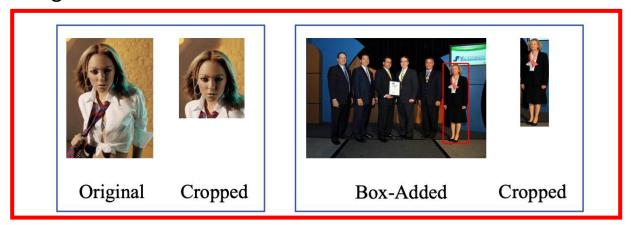


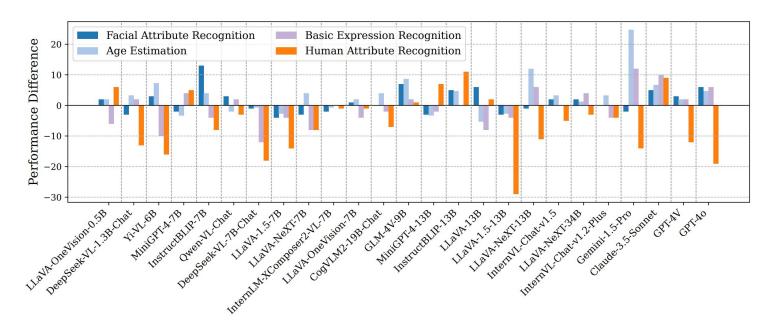






Relative Position of Targets



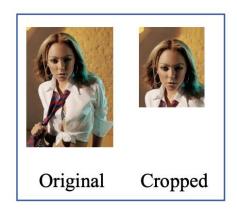




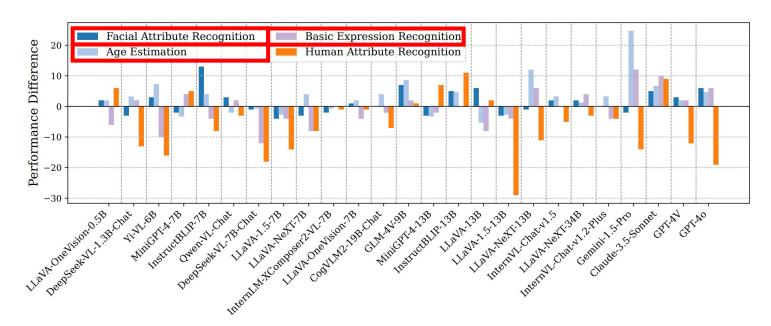




Relative Position of Targets





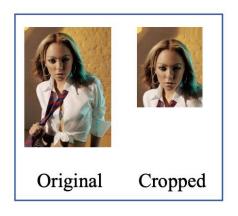




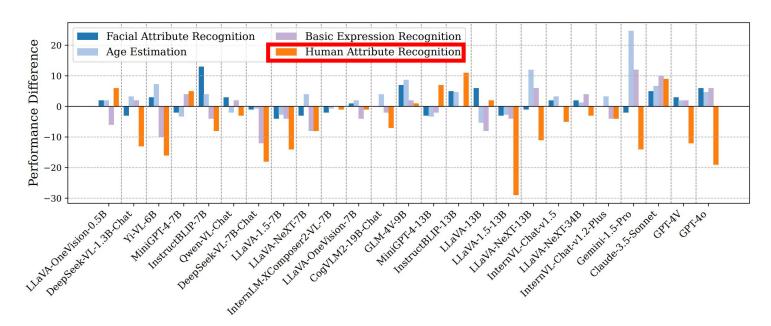




Relative Position of Targets













Setting	Format	Open-	-Source: In	nternVI	ـ-Chat-۱	1.2-Plus		Close-S	ource:	GPT-40	Rea. Overall 71.7 70.0 78.0 73.4 77.2 78.6 81.2 79.9		
Setting	rollilat	Face	Human	Per.	Rea.	Overall	Face	Human	Per.	Rea.	Overall		
ZS	QO→A	69.7	83.1	76.7	76.0	76.4	68.5	71.6	68.9	71.7	70.0		
H	QOH→A	68.4	83.2	76.4	75.9	75.9	72.2	74.6	70.4	78.0	73.4		
H+VCoT	QOH→RA	69.1	82.5	75.9	74.8	75.7	76.4	80.7	78.2	77.2	78.6		
H+1TCoT	QOH→RA	68.6	81.4	75.6	74.3	75.0	77.9	81.9	79.0	81.2	79.9		
H+2TCoT	$QOH \rightarrow R, QOHR \rightarrow A$	69.1	79.1	75.8	71.8	74.1	77.0	81.2	78.4	77.2	79.1		







Satting Farmet		Open-Source: InternVL-Chat-v1.2-Plus					Close-Source: GPT-4o				
Setting	Format	Face	Human	Per.	Rea.	Overall	Face	Human	Per.	Rea.	Overall
ZS	QO→A	69.7	83.1	76.7	76.0	76.4	68.5	71.6	68.9	71.7	70.0
H	QOH→A	68.4	83.2	76.4	75.9	75.9	72.2	74.6	70.4	78.0	73.4
H+VCoT	QOH→RA	69.1	82.5	75.9	74.8	75.7	76.4	80.7	78.2	77.2	78.6
H+1TCoT	QOH→RA	68.6	81.4	75.6	74.3	75.0	77.9	81.9	79.0	81.2	79.9
H+2TCoT	$QOH \rightarrow R, QOHR \rightarrow A$	69.1	79.1	75.8	71.8	74.1	77.0	81.2	78.4	77.2	79.1







Catting	Format	Open-	Source: Ir	iternVL	-Chat-v	1.2-Plus		Close-S	ource:	GPT-40	ÿ.
Setting	Format	Face	Human	Per.	Rea.	Overall	Face	Human	Per.	Rea.	Overall
ZS	QO→A	69.7	83.1	76.7	76.0	76.4	68.5	71.6	68.9	71.7	70.0
H	QOH→A	68.4	83.2	76.4	75.9	75.9	72.2	74.6	70.4	78.0	73.4
H+VCoT	QOH→RA	69.1	82.5	75.9	74.8	75.7	76.4	80.7	78.2	77.2	78.6
H+1TCoT	QOH→RA	68.6	81.4	75.6	74.3	75.0	77.9	81.9	79.0	81.2	79.9
H+2TCoT	$QOH \rightarrow R, QOHR \rightarrow A$	69.1	79.1	75.8	71.8	74.1	77.0	81.2	78.4	77.2	79.1
ale S		à									



Answering Q2: Specialist Models vs. MLLMs





L3 Ability	Age	Exp	ression	Deepfake	Spoofing	Action	Counting
Dataset	UTKFace	RAF-DB (Basic)	RAF-DB (Compound)	FF++	SiW-Mv2	HICO-DET	ShTech-A
Matric	MAE↓	ACC ↑	ACC ↑	ACC ↑	$ACER \downarrow$	mAP ↑	$MAE \downarrow$
Random	27.89	13.85	8.08	50.84	50.05	9.32	1512.65
InternVL-Chat-v1.5	6.43	72.23	42.93	56.21	14.84	22.29	2195.69
LLaVA-NeXT-34B	6.01	77.71	41.04	53.42	22.38	13.74	1592.55
InternVL-Chat-v1.2-Plus	5.21	76.40	30.56	52.89	19.92	12.25	2518.25
Best of The Above 3	5.21	77.71	42.93	56.21	14.84	22.29	1592.55
Early Specialist Model	5.47	74.20	44.55	82.01	9.40	19.81	110.20
Relative Score	1.01	1.06	0.96	0.17	0.87	1.24	-0.06
Need Specialist Model?	No.	No.	No.	Yes.	No.	No.	Yes.
L3 Ability	Basic FR	C.P. FR	C.A. FR	S.L. FR	Occ. FR		Re-ID
Dataset	LFW	CPLFW	CALFW	SLLFW	MLFW		Market1501
Matric	ACC ↑	ACC ↑	ACC ↑	ACC ↑	ACC ↑		$ACC \uparrow^{12}$
Random	50.05	49.75	50.12	50.18	50.05		49.47
InternVL-Chat-v1.5	83.68	58.13	61.40	56.72	52.15		77.53
LLaVA-NeXT-34B	91.32	65.87	62.07	70.25	53.73		85.67
InternVL-Chat-v1.2-Plus	92.57	67.98	66.50	68.50	58.65		88.73
Best of The Above 3	92.57	67.98	66.50	70.25	58.65		88.73
Early Specialist Model	99.50	87.47	92.43	98.40	82.87		95.26
Relative Score	0.86	0.48	0.39	0.42	0.26		0.86
Need Specialist Model?	No.	Yes.	Yes.	Yes.	Yes.		No.



Answering Q2: Specialist Models vs. MLLMs





L3 Ability	Age		ression	Deepfake	Spoofing	Action	Counting
Dataset	UTKFace	RAF-DB (Basic)	RAF-DB (Compound)	FF++	SiW-Mv2	HICO-DET	ShTech-A
Matric	MAE↓	ACC ↑	ACC ↑	ACC ↑	$ACER \downarrow$	mAP ↑	$MAE \downarrow$
Random	27.89	13.85	8.08	50.84	50.05	9.32	1512.65
InternVL-Chat-v1.5	6.43	72.23	42.93	56.21	14.84	22.29	2195.69
LLaVA-NeXT-34B	6.01	77.71	41.04	53.42	22.38	13.74	1592.55
InternVL-Chat-v1.2-Plus	5.21	76.40	30.56	52.89	19.92	12.25	2518.25
Best of The Above 3	5.21	77.71	42.93	56.21	14.84	22.29	1592.55
Early Specialist Model	5.47	74.20	44.55	82.01	9.40	19.81	110.20
Relative Score	1.01	1.06	0.96	0.17	0.87	1.24	-0.06
Need Specialist Model?	No.	No.	No.	Yes.	No.	No.	Yes.
L3 Ability	Basic FR	C.P. FR	C.A. FR	S.L. FR	Occ. FR		Re-ID
Dataset	LFW	CPLFW	CALFW	SLLFW	MLFW		Market1501
Matric	ACC ↑	ACC ↑	ACC ↑	ACC ↑	ACC ↑		ACC \uparrow^{12}
Random	50.05	49.75	50.12	50.18	50.05		49.47
InternVL-Chat-v1.5	83.68	58.13	61.40	56.72	52.15		77.53
LLaVA-NeXT-34B	91.32	65.87	62.07	70.25	53.73		85.67
InternVL-Chat-v1.2-Plus	92.57	67.98	66.50	68.50	58.65		88.73
Best of The Above 3	92.57	67.98	66.50	70.25	58.65		88.73
Early Specialist Model	99.50	87.47	92.43	98.40	82.87		95.26
Relative Score	0.86	0.48	0.39	0.42	0.26		0.86
Need Specialist Model?	No.	Yes.	Yes.	Yes.	Yes.		No.



Answering Q2: Specialist Models vs. MLLMs





L3 Ability	Age		ression	Deepfake	Spoofing	Action	Counting
Dataset	UTKFace	RAF-DB (Basic)	RAF-DB (Compound)	FF++	SiW-Mv2	HICO-DET	ShTech-A
Matric	MAE↓	ACC ↑	ACC ↑	ACC ↑	ACER↓	mAP ↑	$MAE \downarrow$
Random	27.89	13.85	8.08	50.84	50.05	9.32	1512.65
InternVL-Chat-v1.5	6.43	72.23	42.93	56.21	14.84	22.29	2195.69
LLaVA-NeXT-34B	6.01	77.71	41.04	53.42	22.38	13.74	1592.55
InternVL-Chat-v1.2-Plus	5.21	76.40	30.56	52.89	19.92	12.25	2518.25
Best of The Above 3	5.21	77.71	42.93	56.21	14.84	22.29	1592.55
Early Specialist Model	5.47	74.20	44.55	82.01	9.40	19.81	110.20
Relative Score	1.01	1.06	0.96	0.17	0.87	1.24	-0.06
Need Specialist Model?	No.	No.	No.	Yes.	No.	No.	Yes.
L3 Ability	Basic FR	C.P. FR	C.A. FR	S.L. FR	Occ. FR		Re-ID
Dataset	LFW	CPLFW	CALFW	SLLFW	MLFW		Market1501
Matric	ACC ↑	ACC ↑	ACC ↑	ACC ↑	ACC ↑		ACC \uparrow^{12}
Random	50.05	49.75	50.12	50.18	50.05		49.47
InternVL-Chat-v1.5	83.68	58.13	61.40	56.72	52.15		77.53
LLaVA-NeXT-34B	91.32	65.87	62.07	70.25	53.73		85.67
InternVL-Chat-v1.2-Plus	92.57	67.98	66.50	68.50	58.65		88.73
Best of The Above 3	92.57	67.98	66.50	70.25	58.65		88.73
Early Specialist Model	99.50	87.47	92.43	98.40	82.87		95.26
Relative Score	0.86	0.48	0.39	0.42	0.26		0.86
Need Specialist Model?	No.	Yes.	Yes.	Yes.	Yes.		No.









Thanks!

Project Homepage: https://face-human-bench.github.io/ You can also scan the QR code to visit it!



