

Jailbreaking Ensembles of $N = 8$ VLMs

Evaluated VLM

DeepSeek-VL Chat 7B + SigLIP&SAM-B			0.00	
Gemma Instr 2B + CLIP	0.55	0.82	0.78	0.26
Gemma Instr 2B + DINOv2&SigLIP	0.61	0.22	0.26	0.76
Gemma Instr 2B + SigLIP	0.31	0.38	0.25	0.32
Gemma Instr 8B + CLIP	0.80	0.69	0.66	0.80
Gemma Instr 8B + SigLIP	0.20	0.66	0.22	0.21
LLAVAv1.5 7B + CLIP (Repro)	0.32	0.33	0.70	0.49
LLAVAv1.513B + CLIP (Repro)	0.09	0.15	0.13	0.07
Llama2 13B + CLIP	0.47	0.49	0.52	0.47
Llama2 13B + CLIP (Control)	0.68	0.69	0.72	0.68
Llama2 13B + DINOv2&SigLIP	0.48	0.37	0.34	0.51
Llama2 13B + DINOv2&SigLIP (Control)		0.56	0.55	0.62
Llama2 13B + SigLIP	0.42	0.42	0.40	0.48
Llama2 13B + SigLIP (Control)	0.66	0.69	0.63	0.80
Llama2 7B + CLIP	0.51	0.50	0.23	0.50
Llama2 7B + CLIP (Control)	0.72	0.70	0.85	0.71
Llama2 7B + DINOv2&SigLIP	0.42	0.34	0.32	0.40
Llama2 7B + DINOv2&SigLIP (Control)	0.72	0.67	0.63	0.74
Llama2 7B + SigLIP	0.32	0.38	0.32	0.38
Llama2 7B + SigLIP (Control)	0.66	0.69	0.63	0.85
Llama2 Chat 7B + CLIP	0.16	0.68	0.59	0.14
Llama2 Chat 7B + DINOv2/SigLIP	0.45	0.07	0.07	0.14
Llama2 Chat 7B + SigLIP	0.13	0.25	0.13	0.36
Llama3 Instr 8B + CLIP	0.75	0.27	0.65	0.59
Llama3 Instr 8B + DINOv2/SigLIP	0.32	0.21	0.17	0.29
Llama3 Instr 8B + SigLIP	0.27	0.62	0.27	0.66
Mistral Instr v0.2 7B + CLIP	0.83	0.78	0.79	0.81
Mistral Instr v0.2 7B + DINOv2/SigLIP	0.77	0.48	0.59	0.65
Mistral Instr v0.2 7B + SigLIP	0.59	0.67	0.63	0.64
Qwen VL Chat	0.05	0.05	0.05	0.68



Gemma Instr 2B + CLIP	Gemma Instr 2B + CLIP	Llama2 7B + CLIP	Qwen VL Chat
Gemma Instr 2B + DINOv2&SigLIP	Gemma Instr 2B + SigLIP	Llama2 7B + CLIP (Control)	DeepSeek-VL Chat 7B + SigLIP&SAM-B
Gemma Instr 8B + CLIP	Gemma Instr 8B + CLIP	Gemma Instr 2B + CLIP	Gemma Instr 2B + DINOv2&SigLIP
Gemma Instr 2B + DINOv2&SigLIP	Gemma Instr 8B + SigLIP	Gemma Instr 8B + CLIP	Gemma Instr 8B + CLIP
Llama2 Chat 7B + CLIP	Llama2 Chat 7B + CLIP	Llama2 Chat 7B + CLIP	Gemma Instr 2B + DINOv2&SigLIP
Llama2 Chat 7B + DINOv2/SigLIP	Llama2 Chat 7B + SigLIP	Llama3 Instr 8B + CLIP	Llama2 Chat 7B + SigLIP
Llama3 Instr 8B + CLIP	Llama3 Instr 8B + CLIP	Mistral Instr v0.2 7B + CLIP	Llama3 Instr 8B + SigLIP
Llama3 Instr 8B + DINOv2/SigLIP	Llama3 Instr 8B + SigLIP	LLAVAv1.5 7B + CLIP (Repro)	LLAVAv1.5 7B + CLIP (Repro)

Attacked Ensemble of $N = 8$ VLMs