

Jailbreaking Ensembles of $N = 8$ VLMs

Evaluated VLM

DeepSeek-VL Chat 7B + SigLIP&SAM-B			1.00	
Gemma Instr 2B + CLIP	0.78	0.83	0.80	0.82
Gemma Instr 2B + DINOv2&SigLIP	0.51	0.84	0.84	0.31
Gemma Instr 2B + DINOv2&SigLIP	0.90	0.81	0.82	0.78
Gemma Instr 2B + SigLIP	0.77	0.74	0.82	0.76
Gemma Instr 8B + CLIP	0.94	0.90	0.89	0.81
Gemma Instr 8B + SigLIP	0.83	0.51	0.82	0.83
LLAVAv1.5 7B + CLIP (Repro)	0.70	0.69	0.29	0.65
LLAVAv1.513B + CLIP (Repro)	0.95	0.91	0.92	0.93
Llama2 13B + CLIP	0.64	0.60	0.58	0.61
Llama2 13B + CLIP (Control)	0.35	0.38	0.34	0.37
Llama2 13B + DINOv2&SigLIP	0.60	0.69	0.74	0.58
Llama2 13B + DINOv2&SigLIP (Control)		0.46	0.47	0.46
Llama2 13B + SigLIP	0.67	0.65	0.69	0.60
Llama2 13B + SigLIP (Control)	0.42	0.40	0.45	0.26
Llama2 7B + CLIP	0.61	0.63	0.94	0.62
Llama2 7B + CLIP (Control)	0.37	0.38	0.23	0.40
Llama2 7B + DINOv2&SigLIP	0.73	0.80	0.79	0.73
Llama2 7B + DINOv2&SigLIP (Control)	0.40	0.41	0.45	0.36
Llama2 7B + SigLIP	0.77	0.75	0.81	0.74
Llama2 7B + SigLIP (Control)	0.43	0.40	0.45	0.22
Llama2 Chat 7B + CLIP	0.93	0.39	0.42	0.87
Llama2 Chat 7B + DINOv2/SigLIP	0.60	0.94	0.95	0.89
Llama2 Chat 7B + SigLIP	0.89	0.79	0.91	0.68
Llama3 Instr 8B + CLIP	0.30	0.85	0.47	0.46
Llama3 Instr 8B + DINOv2/SigLIP	0.76	0.86	0.87	0.74
Llama3 Instr 8B + SigLIP	0.76	0.45	0.75	0.39
Mistral Instr v0.2 7B + CLIP	0.26	0.32	0.30	0.26
Mistral Instr v0.2 7B + DINOv2/SigLIP	0.33	0.57	0.47	0.40
Mistral Instr v0.2 7B + SigLIP	0.46	0.51	0.43	0.47
Qwen VL Chat	0.97	0.98	0.98	0.52



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