

Jailbreaking Ensembles of $N = 8$ VLMs

Evaluated VLM

Gemma Instr 2B + CLIP	3.32	3.39	3.41
Gemma Instr 2B + DINOv2&SigLIP	3.54	3.63	3.60
Gemma Instr 2B + SigLIP	3.63	3.21	3.61
Gemma Instr 8B + CLIP	3.43	3.44	3.58
Gemma Instr 8B + SigLIP	3.67	3.37	3.66
LLAVAv1.5 7B + CLIP (Repro)	2.86	2.88	2.60
LLAVAv1.513B + CLIP (Repro)	2.68	2.71	2.64
Llama2 13B + CLIP	2.69	2.68	2.67
Llama2 13B + CLIP (Control)	2.60	2.60	2.57
Llama2 13B + DINOv2&SigLIP	2.83	2.83	2.81
Llama2 13B + DINOv2&SigLIP (Control)	2.52	2.55	2.54
Llama2 13B + SigLIP	2.79	2.78	2.78
Llama2 13B + SigLIP (Control)	2.64	2.62	2.64
Llama2 7B + CLIP	2.79	2.83	2.91
Llama2 7B + CLIP (Control)	2.61	2.63	2.57
Llama2 7B + DINOv2&SigLIP	2.83	2.86	2.83
Llama2 7B + DINOv2&SigLIP (Control)	2.63	2.67	2.67
Llama2 7B + SigLIP	2.86	2.81	2.84
Llama2 7B + SigLIP (Control)	2.64	2.62	2.64
Llama2 Chat 7B + CLIP	2.74	2.73	2.70
Llama2 Chat 7B + DINOv2/SigLIP	2.71	2.87	2.86
Llama2 Chat 7B + SigLIP	2.86	2.78	2.87
Llama3 Instr 8B + CLIP	3.19	2.95	2.91
Llama3 Instr 8B + DINOv2/SigLIP	2.98	2.91	2.90
Llama3 Instr 8B + SigLIP		2.97	2.96
Mistral Instr v0.2 7B + CLIP	3.16	3.18	3.18
Mistral Instr v0.2 7B + SigLIP		3.15	
Qwen VL Chat	3.61	3.59	3.62

Cross Entropy Loss

Gemma Instr 2B + CLIP
Gemma Instr 2B + DINOv2&SigLIP
Gemma Instr 2B + DINOv2&SigLIP
Gemma Instr 8B + CLIP
Llama2 Chat 7B + CLIP
Llama2 Chat 7B + DINOv2/SigLIP
Llama3 Instr 8B + CLIP
Llama3 Instr 8B + DINOv2/SigLIP

Gemma Instr 2B + CLIP
Gemma Instr 2B + SigLIP
Gemma Instr 8B + CLIP
Gemma Instr 8B + SigLIP
Llama2 Chat 7B + CLIP
Llama2 Chat 7B + SigLIP
Llama3 Instr 8B + CLIP
Llama3 Instr 8B + SigLIP

Gemma Instr 2B + CLIP
Gemma Instr 8B + CLIP
LLAVAv1.5 7B + CLIP (Repro)
Llama2 7B + CLIP
Llama2 7B + CLIP (Control)
Llama2 Chat 7B + CLIP
Llama3 Instr 8B + CLIP
Mistral Instr v0.2 7B + CLIP

Attacked Ensemble of $N = 8$ VLMs