

# Jailbreaking Ensembles of $N = 8$ VLMs

Evaluated VLM

Gemma Instr 2B + CLIP		0.22
Gemma Instr 2B + DINOv2&SigLIP	0.25	0.84
Gemma Instr 2B + DINOv2&SigLIP	0.25	0.82
Gemma Instr 2B + SigLIP	0.76	0.18
Gemma Instr 8B + CLIP	0.28	0.15
Gemma Instr 8B + SigLIP	0.89	0.16
LLAVAv1.5 7B + CLIP (Repro)	0.79	0.77
LLAVAv1.5 13B + CLIP (Repro)	0.94	0.93
Llama2 13B + CLIP	0.59	0.65
Llama2 13B + CLIP (Control)	0.41	0.37
Llama2 13B + DINOv2&SigLIP	0.61	0.71
Llama2 13B + DINOv2&SigLIP (Control)		0.36
Llama2 13B + SigLIP	0.58	
Llama2 13B + SigLIP (Control)	0.45	0.23
Llama2 7B + CLIP	0.67	0.71
Llama2 7B + CLIP (Control)	0.36	0.47
Llama2 7B + DINOv2&SigLIP	0.69	0.73
Llama2 7B + DINOv2&SigLIP (Control)	0.35	0.40
Llama2 7B + SigLIP	0.74	0.54
Llama2 7B + SigLIP (Control)	0.41	0.24
Llama2 Chat 7B + CLIP	0.27	0.22
Llama2 Chat 7B + DINOv2/SigLIP	0.17	0.93
Llama2 Chat 7B + SigLIP	0.87	0.32
Llama3 Instr 8B + CLIP	0.17	
Llama3 Instr 8B + DINOv2/SigLIP	0.16	0.78
Llama3 Instr 8B + SigLIP		0.22
Mistral Instr v0.2 7B + CLIP	0.35	0.31
Mistral Instr v0.2 7B + DINOv2/SigLIP	0.40	0.47
Mistral Instr v0.2 7B + SigLIP	0.47	0.52
Qwen VL Chat	0.97	0.96

Attack Failure Rate

Gemma Instr 2B + CLIP  
Gemma Instr 2B + DINOv2&SigLIP  
Gemma Instr 8B + CLIP  
Gemma Instr 2B + DINOv2&SigLIP  
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

Attacked Ensemble of  $N = 8$  VLMs