Jailbreaking Ensembles of N=8 VLMs DeepSeek-VL Chat 7B + SigLIP&SAM-B 0.07 0.07 Gemma Instr 2B + CLIP 0.07 0.83 0.80 Gemma Instr 2B + DINOv2&SigLIP 0.15 0.71 Gemma Instr 2B + DINOv2&SigLIP 0.08 0.14 Gemma Instr 2B + SigLIP 0.63 0.78 0.09 0.13 0.08 Gemma Instr 8B + CLIP 0.17 0.79 Gemma Instr 8B + SigLIP0.81 0.74 LLAVAv1.5 7B + CLIP (Repro) 0.73 0.11 LLAVAv1.513B + CLIP (Repro) 0.95 0.94 0.89 Llama2 13B + CLIP 0.58 0.61 0.44 Llama2 13B + CLIP

Llama2 13B + CLIP (Control)

Llama2 13B + DINOv2&SigLIP

Llama2 13B + DINOv2&SigLIP (Control) 0.29 0.27 0.17 0.57 0.67 0.65 0.31 0.55 Llama2 13B + SigLIP 0.53 0.55 0.30 0.32 0.37 Llama2 13B + SigLIP (Control) Evaluated 0.52 Llama2 7B + CLIP 0.59 0.08 0.30 0.30 Llama2 7B + CLIP (Control) 0.08 Llama2 7B + DINOv2&SigLIP 0.63 0.62 0.64 0.36 0.26 0.32 Llama2 7B + DINOv2&SigLIP (Control) 0.59 Llama2.7B + SigLIP0.69 0.63 0.30 0.37 Llama2 7B + SigLIP (Control) 0.35 Llama2 Chat 7B + CLIP 0.09 0.10 0.12 Llama2 Chat 7B + DINOv2/SigLIP 0.08 0.93 0.92 Llama2 Chat 7B + SigLIP 0.85 0.12 0.88 0.10 0.12 0.12 Llama3 Instr 8B + CLIP Llama3 Instr 8B + DINOv2/SigLIP 0.86 0.84 Llama3 Instr 8B + SigLIP 0.73 0.71 0.30 Mistral Instr v0.2 7B + CLIP 0.31 0.11 Mistral Instr v0.2 7B + DINOv2/SigLIP 0.31 0.39 0.40 0.36 Mistral Instr v0.2 7B + SigLIP 0.45 0.49 0.97 0.96 Qwen VL Chat Gemma Instr 2B + CLIPGemma Instr 2B + CLIPLlama2 7B + CLIP Gemma Instr 2B + DINOv2&SigLIP Gemma Instr 2B + SigLIPLlama2 7B + CLIP (Control) Gemma Instr 8B + CLIP Gemma Instr 8B + CLIP Gemma Instr 2B + CLIP Gemma Instr 2B + DINOv2&SigLIP Gemma Instr 8B + SigLIPGemma Instr 8B + CLIP Llama2 Chat 7B + SigLIPLlama2 Chat 7B + DINOv2/SigLIP Llama3 Instr 8B + CLIP Llama3 Instr 8B + CLIP Llama3 Instr 8B + CLIP Mistral Instr v0.2 7B + CLIP Llama3 Instr 8B + SigLIP Llama3 Instr 8B + DINOv2/SigLIP LLAVAv1.5 7B + CLIP (Repro)

Attacked Ensemble of N=8 VLMs

-0.8 -0.6

-0.2