Jailbreaking Ensembles of N=8 VLMs Gemma Instr 2B + CLIP0.19 0.19 0.20 Gemma Instr 2B + SigLIP0.30 LLAVAv1.5 7B + CLIP (Repro) 0.84 0.76 0.22 LLAVAv1.513B + CLIP (Repro) 0.88 0.85 0.82 Llama2 13B + CLIP 0.51 0.60 0.43 Evaluated VLM Llama2 13B + CLIP (Control) 0.65 0.62 0.46 Llama2 13B + DINOv2&SigLIP 0.67 0.60 0.55 Llama2 13B + DINOv2&SigLIP (Control) 0.57 0.59 0.64 Llama $2\ 13B + SigLIP$ 0.69 0.71 0.71 Llama2 13B + SigLIP (Control) 0.47 0.46 0.52 Llama2 7B + CLIP0.58 0.66 0.17 Llama2 7B + CLIP (Control) 0.49 0.16 0.44 Llama2.7B + DINOv2&SigLIP0.62 0.76 0.76 Llama2 7B + DINOv2&SigLIP (Control) 0.47 0.55 0.41 Llama2.7B + SigLIP0.57 0.77 0.77 Llama2 7B + SigLIP (Control) 0.47 0.46 0.52 0.99 0.99 Qwen VL Chat 1.00 Gemma Instr 2B + CLIPGemma Instr 2B + CLIPLlama2 7B + CLIP Gemma Instr 2B + DINOv2&SigLIP Gemma Instr 2B + SigLIPLlama2 7B + CLIP (Control) Gemma Instr 8B + CLIPGemma Instr 8B + CLIP Gemma Instr 2B + CLIP ${\sf Gemma\ Instr\ 2B+DINOv2\&SigLIP}$ Gemma Instr 8B + SigLIPGemma Instr 8B + CLIP Llama2 Chat 7B + SigLIP Llama2 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 $\overset{0}{\overset{0}{7}}$ $\overset{0}{\overset{0}{7}}$ Attack Failure Rate -0.2